



EBERHARD KARLS UNIVERSITÄT TÜBINGEN

MATHEMATISCH-NATURWISSENSCHAFTLICHE FAKULTÄT

WILHELM-SCHICKARD-INSTITUT FÜR INFORMATIK

BACHELORARBEIT INFORMATIK

# How-Provenance Through Query Rewriting

*Pascal Engel*

August 14, 2020

**Gutachter**

Prof. Dr. Torsten GRUST

MATHEMATISCH-NATURWISSENSCHAFTLICHE FAKULTÄT

WILHELM-SCHICKARD-INSTITUT FÜR INFORMATIK

**Betreuer**

Dr. Tobias MÜLLER

MATHEMATISCH-NATURWISSENSCHAFTLICHE FAKULTÄT

WILHELM-SCHICKARD-INSTITUT FÜR INFORMATIK

**Pascal Engel**

*How-Provenance Through Query Rewriting*

Bachelorarbeit Informatik

Eberhard Karls Universität Tübingen

Matrikelnummer: 4118448

Bearbeitungszeit: July 2020 – November 2020

## **Abstract**

For SQL queries in databases, how-provenance can add valuable insight into which expressions of the query are responsible for which cell values from the result. In this work, a two-phase approach developed by the Database Research Group of the University of Tübingen is used to translate a given query into two distinct queries that compute its how-provenance. This work presents a Haskell program to automate this rewriting process for a multitude of SQL constructs.



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Provenance in Databases</b>	<b>3</b>
2.1	SQL . . . . .	3
2.1.1	Queries . . . . .	3
2.1.2	User Defined Functions . . . . .	4
2.1.3	PostgreSQL . . . . .	4
2.2	Data Provenance . . . . .	5
2.2.1	How-Provenance . . . . .	5
2.2.2	Phases of Provenance . . . . .	5
2.3	Example . . . . .	7
<b>3</b>	<b>Formalization</b>	<b>9</b>
3.1	Notation . . . . .	9
3.2	Prerequisites . . . . .	10
3.2.1	Provenance Relations . . . . .	10
3.2.2	Normalized Queries . . . . .	11
3.2.3	Annotations . . . . .	12
3.2.4	Logging . . . . .	12
3.3	Rules . . . . .	12
3.3.1	Operators . . . . .	13
3.3.2	Case . . . . .	14
3.3.3	SFW: $n$ -way join . . . . .	15
<b>4</b>	<b>Implementation</b>	<b>17</b>
4.1	Haskell . . . . .	17

4.1.1	Abstract Syntax Tree . . . . .	18
4.2	Project Structure . . . . .	18
4.2.1	Modules . . . . .	18
4.2.2	Annotate . . . . .	20
4.2.3	Provenance Relations . . . . .	20
4.2.4	Correlated Subqueries . . . . .	20
4.2.5	Provenance Sets . . . . .	21
4.2.6	Push . . . . .	23
4.3	Translation . . . . .	26
4.3.1	Operators . . . . .	27
4.3.2	Case . . . . .	28
4.3.3	<i>n</i> -way-Join . . . . .	29
<b>5</b>	<b>Outlook &amp; Conclusion</b>	<b>33</b>
	<b>Bibliography</b>	<b>35</b>
	<b>Appendices</b>	
<b>A</b>	<b>List of Translation Rules</b>	<b>39</b>
A.1	Literals . . . . .	39
A.2	Operators . . . . .	39
A.3	Column References . . . . .	39
A.4	Case . . . . .	40
A.5	Row Constructors . . . . .	40
A.6	Values . . . . .	41
A.7	Union All . . . . .	41
A.8	With . . . . .	41
A.9	SFW: Plain Select From . . . . .	42
A.10	SFW: <i>n</i> -way join . . . . .	42
A.11	SFW: Aggregates . . . . .	43
A.12	SFW: Order By . . . . .	44
<b>B</b>	<b>Sample Output</b>	<b>45</b>

# List of Figures

- 2.1 Computation of How-Provance using 2 phases . . . . . 6
- 3.1 Phase1/2 relations for table flights . . . . . 11
- 4.1 Modular dependencies . . . . . 19





# 1 | Introduction

With the ever growing amount of data provided today through the internet, it is increasingly important in databases to ensure their trustworthiness and integrity. Here, provenance may help to add value to data by showing its source. [CCT09] Specifically, the notion of *how*-provenance may help debugging a query and finding flaws that lead to unexpected behaviour. Moreover, it is useful to simply understand what query expressions do and how they contribute to a specific result value.

Listing 1.1 shows a *SQL* query that selects the destination of all excessively cheap and expensive inner-European flights from a database (Table 1.1). The query is composed of two subqueries that are being unified. Usually, when running the query we will only be provided with the result table. What *how*-provenance offers here is insight into which part of the query is responsible for which result value. For the sake of simplicity in this example, we are only concerned with subqueries. The first subquery colored in **orange** selects those destinations of flights with a low price. Analogously, the result values that went through the second subquery with a high price are colored **green**.

```
1 (SELECT destination
2 FROM flights
3 WHERE price < 100)
4 UNION ALL
5 (SELECT destination
6 FROM flights
7 WHERE price > 200);
```

**Listing 1.1:** Example query

In Table 1.1 the corresponding values, that are being calculated by it, are colored accordingly. We can see here, exactly which subquery is responsible for which

---

resulting cell value, which is more than *SQL* can provide by itself.

flights				output
origin	destination	distance	price	destination
Porto	Amsterdam	1612	212	Amsterdam
Amsterdam	Munich	668	164	Amsterdam
London	Amsterdam	357	92	Amsterdam
London	Rome	1435	283	Rome
Munich	Rome	698	194	

**Table 1.1:** Input and output table

The example above is only concerned with subquery expressions, but it shows how how-provenance adds insight into the computations of a query, which a Database Management System cannot provide. Using this notion for a variety of additional *SQL* expressions, we can compute more fine-grained provenance, as well.

The goal of this work is to implement a SQL-to-SQL compiler that automates the computation of How-Provenance for any given *SQL* query. For this, the Database Research Group from the University of Tübingen has developed a two-phase approach [OMG18] which lays the theoretical groundwork for this thesis.

In the following chapters, I will begin by explaining Provenance in the context of Databases and the query language *SQL*. Then, the theoretical groundwork, including the main prerequisites and the formalization of the implemented translation rules are presented. Finally, in the main part, the implementation in Haskell is presented and discussed, along with possible further work based on the program.

## 2 | Provenance in Databases

In this chapter, I will outline the Provenance aimed to compute in this work in the context of databases and the associated query language *SQL*.

### 2.1 SQL

The *Structural Query Language (SQL)* has been the dominant language for managing data in databases. It is built on top of the theoretical works on relational algebra. Its syntax is supposed to be intuitive in plain English, so a *SQL* query can be roughly interpreted as an English sentence, making it easily understandable even for people who are not familiar with databases. [CB74]

Besides queries, *SQL* offers various functionalities to create and manipulate tables in a database, as well as guaranteeing their integrity. With the recent addition of recursion it is now even seen as a very capable Turing-complete programming language.

#### 2.1.1 Queries

In *SQL* queries are the essential tool to answer questions about the data present in a given database. The standard form of a query is a *Select From Where (SFW)* expression. In this work a limited *SQL* dialect is used, so not all common *SQL* expressions are covered. *SFW* expressions are limited to the following form:

```
SELECT  $e_1, \dots$   
FROM  $et_1, \dots$   
WHERE  $w$   
GROUP BY  $g_1, \dots$   
ORDER BY  $o_1$   
DISTINCT ON  $d_1, \dots$   
LIMIT  $l_1$   
OFFSET  $l_o$ 
```

In such a query all clauses are optional except the **SELECT** clause which is responsible for returning a result value and thus cannot be left out.

### 2.1.2 User Defined Functions

Generally, *User Defined Functions* (*UDFs*) can be used to compute a result value or manipulate a database. In the scope of this work, they are used without any side-effects. The syntax of these *UDFs* is briefly explained in Listing 2.1, where each  $\tau_i$  stands for a type and  $q$  represents a *SQL* query expression. *UDFs* are more powerful than this template suggests, however in the scope of this work *UDF* definitions of this form are sufficient. Note that here the function arguments are unnamed, so in the function body they will appear in the order they are bound as  $\$1, \$2, \dots$ .

```
1 CREATE [OR REPLACE] FUNCTION  $f(\tau_1, \dots, \tau_n)$  RETURNS  $\tau$  AS  
2 $$  
3    $q$   
4 $$ LANGUAGE SQL;
```

**Listing 2.1:** Syntax for UDF definitions

### 2.1.3 PostgreSQL

PostgreSQL (or simply Postgres) is an open-source relational database management system (DBMS) first introduced in 1996. [pos] Although the query-rewriting is designed independent of the DBMS, Postgres offers some useful functionalities. It generally supports all common *SQL* constructs together with some additions

thanks to its interface for extensions. For these reasons Postgres is the chosen DBMS at the Database Research Group. Here, it is used to run the generated queries from the compiler presented in this work. Additionally, its internal parser is used to parse a query into an Abstract Syntax Tree.

## 2.2 Data Provenance

In general, 'provenance information describes the origin and history of data in its life cycle.' [CCT09, p.1] Provenance can be applied to all kinds of data, wherever it is useful to know how it is created. 'In databases it aims to answer record-level questions, e.g., which tuples (rows) in the input tables contributed to a particular output tuple and how.' [Lud03, p.10] These questions further divide provenance analysis into subcategories concerned with *why* a specific result is being computed, *where* the resulting data originally comes from and *how* it has been computed by the query.

### 2.2.1 How-Provenance

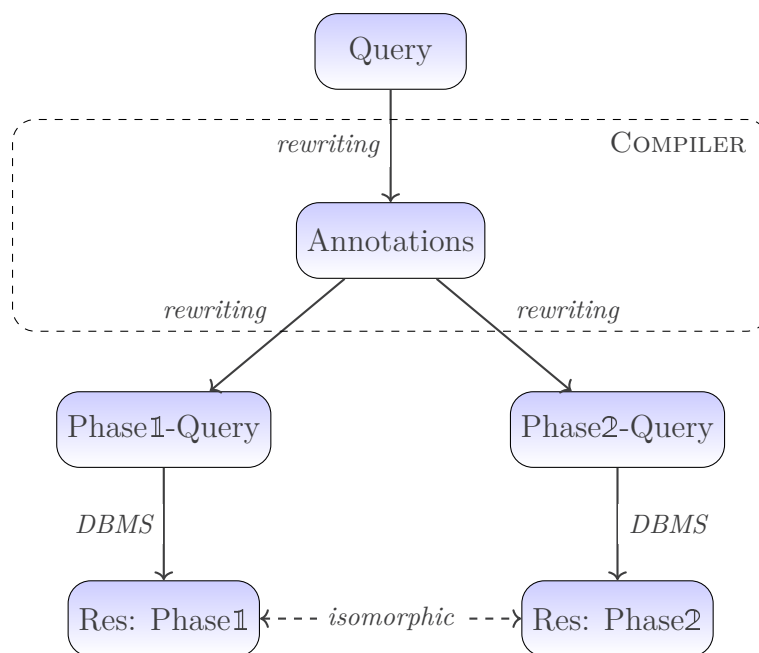
In this work, I focus on the *how*-aspect. The notion of *how*-provenance has first been introduced by *Green et al.* to add insights to data provenance where *why*-provenance could not answer questions sufficiently. [GKT07] *How*-provenance helps tracing output items back through the query and can explain complex queries that build on advanced *SQL* constructs. [OMG18] So analyzing the provenance of a result value will point us at exactly those query constructs that are involved in its computation. It is therefore particularly helpful when debugging queries and generally handy while learning the *SQL* language to gain further insights that a simple result table cannot offer.

### 2.2.2 Phases of Provenance

This work relies on a two-phase approach to translate an input query into two distinct queries, that compute its *how*-provenance. Using this query-rewriting, a database management system can run the generated queries and return the desired provenance data. This is useful because no additional software is required in this

step. A DBMS is already capable of running queries efficiently, so using it for this purpose saves a lot of work.

As shown in Fig. 2.1, in this process each expression of an input query is annotated with a label and translated independently in two phases: In Phase1 the query is being evaluated and as a side-effect corresponding data is logged in a dedicated table. In Phase2 literal data is ignored, instead the query-annotations are *pushed* into the query result while the logging data from Phase1 is being read to ensure correctness. The resulting relations of both phases are tables that are isomorphic to each other (i.e. they consist of the same columns and the same number of rows). This means that for every cell in the result of Phase1, that contains the data of the original query, there is a corresponding cell in the result of Phase2, that contains the collected provenance annotations. Thus, by comparing the annotated query with the resulting relations, we can deduce exactly which part of the query is involved in which resulting cell value.



**Figure 2.1:** Computation of How-Provence using 2 phases

The compiler developed for this thesis aims to automate the rewriting process for both phases. As demonstrated in Fig. 2.1 the compiler covers all rewriting processes: For an input query, an annotated query, a Phase1 query and a Phase2

query are generated. The queries of both phases can then be executed by a DBMS and the resulting tables are isomorphic to each other, so for each cell a provenance set is calculated.

## 2.3 Example

To explain this with a more in-depth example, let's take another look at the query in Listing 1.1. For this rewriting process, the query is first annotated, so each expression has a unique label between 1 and 13 (Listing 2.2).

```

1  2(SELECT 3destination
2    FROM 4flights
3    WHERE 5price 7< 6100)
4  1UNION ALL
5    8(SELECT 9destination
6      FROM 10flights
7      WHERE 11price 13> 12200);

```

Listing 2.2: Annotated query

This query is then transformed based on rules for each expression. The Phase1 translation produces a table containing tuple identifiers and original cell values. For Phase2 the result contains the same tuple identifiers, but the cell values are replaced with all provenance annotations that were involved in the computation for that cell. These annotations are collected in an array for each cell.

result1		result2	
tuid	destination	tuid	destination
3	Amsterdam	3	{1, 2, 3, 4, 5, 6, 7}
1	Amsterdam	1	{1, 8, 9, 10, 11, 12, 13}
5	Rome	5	{1, 8, 9, 10, 11, 12, 13}

Table 2.1: Results of generated queries

We can interpret these two result tables together with the annotated query: For each cell, by inspecting the collected provenance annotations and comparing them to the labels in the annotated query, we can deduce exactly which expressions are involved in its computation. For example, for the cell `Rome` in Phase1, we can compare its `tuid` value with the ones in Phase2 and thus associate the provenance

### 2.3. EXAMPLE

---

set  $\{1, 8, 9, 10, 11, 12, 13\}$  with it. By comparing the elements in the provenance set with the labels in the annotated query, we can deduce that the **UNION ALL** expression and all expressions contained in the second subquery are involved in its computation.

Additionally, for both *Amsterdam* cells in Phase1 we gain the same insight shown in the introduction: The cell with the tuple identifier 3 was computed by the first subquery and the one with 1 by the second subquery.



## 3 | Formalization

In this chapter, the theoretical background for the compiler presented in this work is outlined.

### 3.1 Notation

This section provides an overview of the notations used in the formalization.

**Meta Variables** These meta variables represent different kinds of *SQL* expressions.

Symbol	Description
$e$	Meta variable for an arbitrary expression
$l$	Meta variable for a literal expression
$ec$	Meta variable for a cell expression
$et$	Meta variable for a table expression

**Provenance Data** The data describing provenance is notated as sets. Annotations form the elements of these sets.

Symbol	Description
$\alpha$	Meta variable for an annotation (represented as integer)
$e^\alpha$	Annotated expression
$\mathfrak{p}$	Meta variable for a provenance set
$\{\dots\}$	Set of provenance
$\cup$	Binary set union
$\Psi(e, \mathfrak{p})$	Add data provenance $\mathfrak{p}$ to each provenance set found in $e$ 's result value

**Logging** Various logging functions are used to keep track of data (e.g. while filtering).

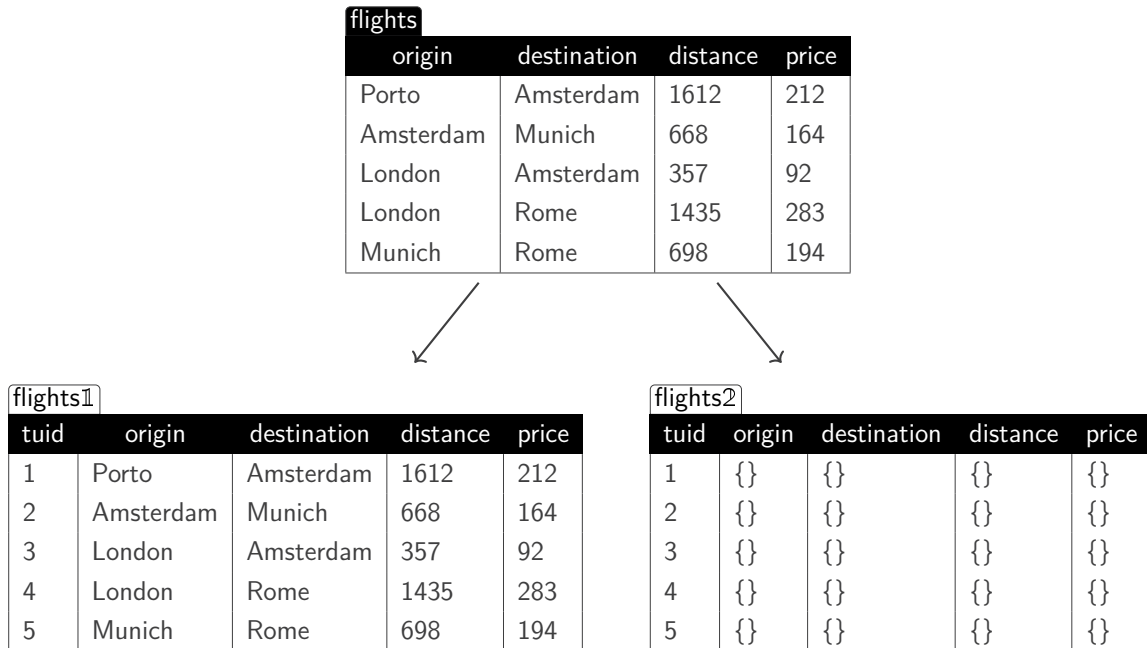
Symbol	Description
$f$	Meta variable for a free tuple variable
$\{f_1, \dots\} := \text{fv}(e)$	Function to find free tuple variables in an expression $e$
$\textcircled{D}$	Text location identifier (to keep track of logging call)
$\text{cast}(\alpha)$	Convert an annotation to a location identifier
$\text{writeX}$	<i>UDF</i> to log data in Phase1 in the context of $x$
$\text{readX}$	<i>UDF</i> to read logs in Phase2 in the context of $x$

## 3.2 Prerequisites

In order for the generated *SQL* code to run properly, a few conditions have to hold, before the query rewriting can begin. These are shown in this section.

### 3.2.1 Provenance Relations

For each table referenced in the query, there have to be two relations concerned with its provenance. These relations replace the original relations referenced in the input query. In this work, they are defined by *SQL* Views. As shown in Fig. 3.1, for the **flights** table in the introduction, the Phase1 relation contains the same data with an additional **tuid** column. The Phase2 relation is made of the same columns and contains the same values as **tuid**, however each original data cell is replaced with an empty provenance set. These provenance sets are finally used to accumulate all relevant annotations for the result.

Figure 3.1: Phase1/2 relations for table `flights`

### 3.2.2 Normalized Queries

SFW-blocks have to be normalized so each SFW-block falls under just one of four categories. SFW-blocks can be a plain select, n-way joins, aggregates and order-by. Queries, that don't hold this, need to be transformed to separate the responsibilities to subqueries. This clause isolation imposes clear advantages here: Instead of defining rules for all possible ways of defining a SFW-block, we only use rules for specific cases of blocks and limit the total number of translation rules and thus achieve a clear separation of responsibilities. Furthermore, this adds expressiveness to the query constructs which may help while debugging.

As described in *Müller et al.*, a given query can be normalized to a query containing multiple nested subqueries with distinct cases. [MDG18, p.1541]

$$\begin{array}{l}
\text{SELECT } e_1, \dots \\
\text{FROM } et_1, \dots \\
\text{WHERE } w \\
\text{GROUP BY } g_1, \dots \\
\text{ORDER BY } o_1 \\
\text{DISTINCT ON } d_1, \dots \\
\text{LIMIT } l_l \\
\text{OFFSET } l_o
\end{array}
\Rightarrow^{norm}
\begin{array}{l}
\text{SELECT } e_1, \dots \\
\text{FROM } ( \\
\quad \text{SELECT } e_1, \dots \\
\quad \text{FROM } ( \\
\quad \quad \text{SELECT } e_1, \dots \\
\quad \quad \text{FROM } et_1, \dots \\
\quad \quad \text{WHERE } w \\
\quad \text{GROUP BY } g_1, \dots \\
\quad ) \\
\text{ORDER BY } o_1, \dots \\
\text{DISTINCT ON } d_1, \dots \\
\text{LIMIT } l_l \\
\text{OFFSET } l_o
\end{array}$$

### 3.2.3 Annotations

Before rewriting can begin, each expression  $e$  needs to be annotated with a label  $\alpha$  to  $e^\alpha$  that represents the expression's provenance data.

### 3.2.4 Logging

For case-expressions, joining, filtering, aggregating and unifying queries, tuples need to be logged. Therefore, dedicated logging-UDFs have to be loaded into the database before running the generated queries. These UDFs write in Phase1 relevant tuple identifiers into a logging table and read them in Phase2, so that branches and filters still apply even when the relevant data is not present. [MDG18, p.1543]

## 3.3 Rules

This work includes the rewriting of most *SFW* expressions, basic operations like function applications, column references, **ROW** and **VALUES** constructors, **CASE** ex-

pressions, **WITH** bindings and **UNION ALL**. The following presents a selection of the provided translation rules in order to exemplary explain the rewriting-process for both phases. A complete list of the translation rules is added in the Appendix. The section above the bar defines the environment under which the rule is applied. This also includes prerequisites for the rule, recursively already translated subexpressions and other bindings for various terms. Under the bar, on the left hand side is a generalized annotated expression as input and the rewritten expression on the right hand side as output expression. A translation rule is denoted as

$$SQL^\alpha \Rightarrow^\tau (SQL^{\mathbb{1}}, SQL^{\mathbb{2}})$$

where  $SQL^\alpha$  is an annotated  $SQL$  expression,  $\tau$  is the type of an  $SQL$  expression and  $(SQL^{\mathbb{1}}, SQL^{\mathbb{2}})$  is the tuple consisting of the translated Phase $\mathbb{1}$  and Phase $\mathbb{2}$   $SQL$  expressions.

The overall rewriting process is the composition of these rules.

### 3.3.1 Operators

This rule is concerned with binary operators on a cell level. These can be arithmetic ( $+$ ,  $*$ ,  $\dots$ ) or logic ( $\&$ ,  $\dots$ ). In Phase $\mathbb{1}$  the generalized operator  $\otimes$  is simply applied to its arguments  $ec_1^{\mathbb{1}}, ec_2^{\mathbb{1}}$ . In Phase $\mathbb{2}$ , because the subexpressions  $ec_1^{\mathbb{2}}, ec_2^{\mathbb{2}}$  are expected to be provenance sets, they are unified to a singular set, where the label of the operator expression is added to. The provenance of the operator is computed by replacing  $\otimes$  with  $\cup\{\alpha\}\cup$ , so instead of the operator its label is added and the resulting how-provenance consists of the provenance of the subexpressions and the label of the operator expression.

$$\frac{\text{BINOP SCALAR} \quad ec_1^{\alpha_1} \Rightarrow^{cell} (ec_1^{\mathbb{1}}, ec_1^{\mathbb{2}}) \quad ec_2^{\alpha_2} \Rightarrow^{cell} (ec_2^{\mathbb{1}}, ec_2^{\mathbb{2}})}{(ec_1^{\alpha_1} \otimes ec_2^{\alpha_2})^\alpha \Rightarrow^{cell} (ec_1^{\mathbb{1}} \otimes ec_2^{\mathbb{1}}, ec_1^{\mathbb{2}} \cup \{\alpha\} \cup ec_2^{\mathbb{2}})}$$

#### 3.3.2 Case

This rule is concerned with the *SQL* **CASE**-construct. The idea for this rule is to label each branch and memorize the one which is taken in Phase1 and read that information in Phase2 to ensure the same branch is being taken, while collecting all provenance for expressions that need to be evaluated for the branch-decision.

In detail, in Phase1 all **WHEN**-expressions are evaluated, until one evaluates to true. Then, the corresponding integer value will be passed to the function `writeCase` in order to log the branch to be taken. After that, the corresponding **THEN**-expression will be evaluated.

In Phase2 the logged integer to determine the branch is loaded with the `readCase` function. Accordingly, the corresponding provenance of the **THEN**-expression is returned, alongside the provenance of *all* **WHEN**-expressions looked at until this point. The provenance of all previous **WHEN**-expressions has to be considered, because all of them have to be evaluated and rejected until the correct branch is being taken. This is why, if the **ELSE**-branch is taken, all **WHEN**-expressions are in the result. Finally, the label  $\alpha$  of the **CASE**-expression is pushed into the result.<sup>1</sup>

---

<sup>1</sup>Note that free Variables are not considered here, however the rule presented in the Appendix covers free variables.

$$\begin{array}{l}
 \text{CASE} \\
 \left| ec_{w_i}^{\alpha_{w_i}} \Rightarrow^{cell} (ec_{w_i}^1, ec_{w_i}^2) \right|_{i=1\dots} \quad \left| ec_{t_i}^{\alpha_{t_i}} \Rightarrow^{cell} (ec_{t_i}^1, ec_{t_i}^2) \right|_{i=0\dots} \\
 \mathcal{O} := \text{tolocation}(\alpha) \\
 \text{CASE writeCase}(\mathcal{O}, \\
 \quad \text{CASE} \\
 \quad \quad \text{WHEN } ec_{w_1}^1 \text{ THEN } 1 \\
 \quad \quad \vdots \\
 \quad \quad \text{ELSE } 0 \\
 X^1 := \quad \text{END} \\
 \quad ) \\
 \quad \text{WHEN } 1 \text{ THEN } ec_{t_1}^1 \\
 \quad \quad \vdots \\
 \quad \quad \text{ELSE } ec_{t_0}^1 \\
 \quad \text{END} \\
 \text{CASE readCase}(\mathcal{O}) \\
 \text{WHEN } 1 \text{ THEN } \Psi(ec_{t_1}^2, ec_{w_1}^2) \\
 \quad \vdots \\
 X^2 := \text{WHEN } i \text{ THEN } \Psi(ec_{t_i}^2, ec_{w_1}^2 \cup \dots \cup ec_{w_i}^2) \\
 \quad \vdots \\
 \quad \text{ELSE } \Psi(ec_{t_0}^2, ec_{w_1}^2 \cup \dots) \\
 \text{END} \\
 \hline
 e_{in} := \left( \begin{array}{l} \text{CASE} \\ \quad \text{WHEN } ec_{w_1}^{\alpha_{w_1}} \text{ THEN } ec_{t_1}^{\alpha_{t_1}} \\ \quad \quad \vdots \\ \quad \quad \text{ELSE } ec_{t_0}^{\alpha_{t_0}} \\ \text{END} \end{array} \right)^\alpha \Rightarrow^{cell} (X^1, \Psi(X^2, \{\alpha\}))
 \end{array}$$

### 3.3.3 SFW: $n$ -way join

The  $n$ -way Join is a centerpiece in relational algebra for combining multiple relations and filtering results. Here, its generalized pattern is given with  $et_{inp}$ . The **FROM** section may yield an arbitrary number of relations, the **WHERE**-expression is optional (will be interpreted as *true* if missing). No additional clauses are permitted in the *SFW* block.

### 3.3. RULES

---

In Phase1 the `writeJoin` *UDF* logs those tuple identifiers, that have passed the filtering and joining process. Exactly those are also returned in the **SELECT**-clause. As Phase2 cannot handle literal data the information of the Phase1-logging has to be read as a table *log* that is joined with the original relations to only include rows with a tuple identifier, which has been previously logged. This way, the joining and filtering process from Phase1 is reconstructed. Additionally, because here we cannot make use of the **WHERE**-expression, its provenance is pushed into each column reference in the **SELECT**-section, so it will be present in the final provenance set. At last, the label  $\alpha$  is *pushed* to the final result of the expression.

SFW-JOIN

$$\begin{aligned}
 & \mathbb{D} := \text{cast}(\alpha) \quad \{f_1, \dots\} := \text{fv}(et_{inp}) \\
 & \quad \left| \begin{array}{l} ec_i^{\alpha_i} \Rightarrow^{cell} (ec_i^1, ec_i^2) \Big|_{i=0, \dots, n} \\ et_i^{\alpha_i} \Rightarrow^{table} (et_i^1, et_i^2) \Big|_{i=(n+1), \dots, (n+m)} \end{array} \right. \\
 X^1 := & \quad \text{SELECT } writeJoin(\mathbb{D}, var_1.\rho, \dots, f_1.\rho, \dots) \text{ AS } \rho \\
 & \quad ec_1^1 \text{ AS } col_1, \dots \\
 & \quad \text{FROM } et_1^1 \text{ AS } var_1, \dots \\
 & \quad \text{WHERE } ec_0^1 \\
 X^2 := & \quad \text{SELECT } log.\rho \text{ AS } \rho \\
 & \quad \Psi(ec_1^2, ec_0^2) \text{ AS } col_1, \dots \\
 & \quad \text{FROM } et_1^2 \text{ AS } var_1, \dots, \\
 & \quad readJoin(\mathbb{D}, var_1.\rho, \dots, f_1.\rho, \dots) \text{ AS } log \\
 \hline
 et_{inp} := & \left( \begin{array}{l} \text{SELECT } ec_1^{\alpha_1} \text{ AS } col_1, \dots \\ \text{FROM } et_{n+1}^{\alpha_{n+1}} \text{ AS } var_{n+1}, \dots \\ \text{WHERE } ec_0^{\alpha_0} \end{array} \right)^{\alpha} \Rightarrow^{table} (X^1, \Psi(X^2, \{\alpha\}))
 \end{aligned}$$



## 4 | Implementation

In this chapter, the main part of this work is discussed: the implementation of the compiler based on the mostly given formalization in Chapter 3. I will begin by presenting the tools and structures used for the program and then explain the project structure and each module in detail.

### 4.1 Haskell

The chosen programming language for the implementation is Haskell. [Jon03] Haskell is a declarative, purely functional language first introduced in 1990. In the scope of this work its functional aspect is very useful for defining the translation rules presented in Section 3.3 because in a Haskell program mathematical syntax can easily be interpreted as a program. Therefore, the relation between the formal rules and the implementation in Haskell is particularly intuitive.

Furthermore, the static type system allows us to check the correctness of the program on type level while developing, which makes testing a lot easier. The type system additionally helps separate the scope of distinct parts of the program. E.g. in this project, only the Main module is capable of handling Input/Output, so the query translation which is done in different modules cannot interfere or be influenced by side effects on the machine. Here, the modularity of Haskell also helps separate responsibilities of independent computations (like Phase1 and Phase2).

The inherent *Laziness* and the *Algebraic Data Types* work together neatly while recursively traversing the syntax tree of a *SQL* expression. The *lazy* evaluation of subexpressions can save a decent amount of runtime here.

Finally, the usage of language extensions provides a more expressive and declara-

tive way of coding according to specific needs. In this work, the language extensions *ViewPatterns*, *RecordWildCards* and *GADTs* have proven themselves useful.

### 4.1.1 Abstract Syntax Tree

Previous work on the topic has already been implemented in Haskell: Most notably the *LogParser* designed by Denis Hirn [Hir17]. The *LogParser* library uses an expansion for Postgres to extract logs from the internal Postgres *SQL*-parser and uses them to create an *Abstract Syntax Tree* (AST) representation in Haskell. The type of the AST is a *Generalized Algebraic Datatype* (GADT), which means it encapsulates multiple data types (in this case different kinds of *SQL* expressions) in one generalized type for all *SQL* expressions. So the type system of Haskell is still useful to distinguish different types of *SQL* expressions.

In the implementation presented in this work, all transformations work on top of this AST. Because typically records in the AST have a lot more fields than are relevant for our purposes, using *RecordWildCards* is a simple way to ignore those fields and pass the corresponding values through by inexplicitly binding them to variables that are not being manipulated.

Furthermore, the library exports the ***transformM*** function which implements an easy full traversal of an AST while applying a monadic function to all nodes. This is very useful to ensure the totality of a function defined with ***transformM***.

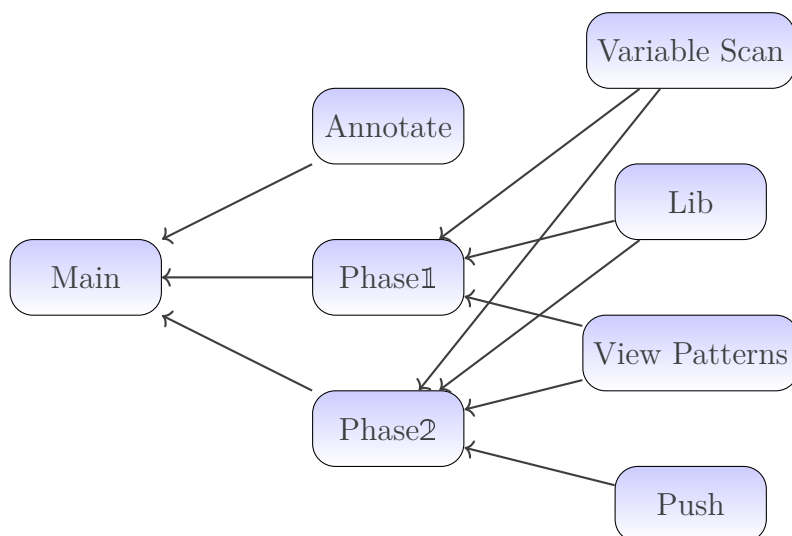
For the output of the program, the *PrettyPrinter* library prints syntax-highlighted *SQL* code from the AST to the console. So it can be checked and inserted into a DBMS to compute provenance.

## 4.2 Project Structure

### 4.2.1 Modules

Fig. 4.1 describes how the structure and internal dependencies of the project presented in this work.

- The *Main* module provides a minimal user interface to load a query, parse it using the *LogParser* and print the annotated query and its *Phase1*- and *Phase2*-translation using the *PrettyPrinter*.



**Figure 4.1:** Modular dependencies

- In the *Annotate* module, the AST is traversed and to each node a globally unique label is added where uniqueness is ensured using a **State** monad as counter.
- In the dedicated modules *Phase1* and *Phase2*, the rules described in Section 3.3 are implemented using plain Haskell pattern-matching as well as the Language Extension *ViewPatterns*.
- The *Lib* module includes all kinds of auxiliary functions, particularly the implementation of provenance sets, tuid-column references and the handling of logging functions.
- The *View Patterns* strictly defines the various forms of SFW-blocks to differentiate in the translation process.
- *Push* includes all functionality associated with the  $\Psi$ -operator and thus enables pushing labels to provenance sets in *Phase2*.
- *VariableScan* is used for getting all free variables in a given subexpression to deal with correlated subqueries.

### 4.2.2 Annotate

In a preprocessing step, labels are distributed to each expression in the AST. *transformM* traverses the AST and each node relevant for *how*-provenance is wrapped in an additional **GLabel**-node. Using a state, it is ensured that every label is unique within the AST.

The labels are only used internally during the translation process. Printed out, the annotated query is not executable because the Label nodes are not part of SQL. However, it needs to be shown to the user in some way in order to interpret the resulting provenance sets.

### 4.2.3 Provenance Relations

As stated in Section 3.2.1, for each table in the query, a relation for both phases needs to be defined. This has already been achieved in related work using materialized views [Par18] and works the same way here: For Phase1 the view consists of the original table with an additional *tuid* column that is filled by a sequence. The Phase2 view references the view from Phase1, so the *tuid* columns have the same values, and replaces all other cells with empty provenance sets. An example for these view definitions is included in the Appendix.

### 4.2.4 Correlated Subqueries

As shown in Section 3.1 the function  $\text{fv}(e)$  is used to get all free tuple variables in an expression  $e$ . This is done to be able to deal with correlated subqueries. A correlated subquery is a subquery that references variables that were outside of its scope, if it was executed separately. In order to deal with these kinds of subqueries, everytime a query-expression is found, it has to be checked, whether there are any free variables present. If so, then all free variables need to be considered while logging, so they need to be passed as arguments to the logging-functions, too.

Finding free variables is done by the function *scan* in the module *VariableScan*, where all definitions of variables and all references of variables are collected and compared in a state monad, so locally undefined tuple variables can be found. This step is done for every SFW-expression (compare Snippets 4.7 and 4.8), which ensures the correctness for correlated subqueries.

```

1 module VariableScan (scan) where
2
3 trav :: Rule (SQL a) (SQL a)
4
5 -- collect var-definition in a table reference
6 trav RTE {...} =
7   modify (<> StateC [aliasname] [])
8   *> transformM trav RTE {...}
9   <* ensure
10
11 -- collect var-reference in a column reference
12 trav EColRef {...} =
13   modify (<> StateC [] [rowVar])
14   *> transformM trav EColRef {...}
15   <* ensure

```

Snippet 4.1: Variable Scan

### 4.2.5 Provenance Sets

Provenance sets are implemented in *SQL* as arrays of integers. This is only an approximation, since arrays do not have the same properties as sets, e.g. they may include duplicate elements. In further work, the implementation type of provenance sets may be subject of change.

On these implemented provenance sets, set operations need to be executed on just like on basic Haskell datatypes, therefore the implementation may look like a **Monoid** instance declaration. However, since the **EArrayExpr** data constructor is only one of many of type **Expr**, a generalized *Monoid* instance can not be defined. Thus, the final implementation as seen in Snippet 4.2 may resemble an instance declaration with the typical functionalities associated with *Monoids*.

```

1 -- ARITHMETICS OF PROVENANCE SETS
2
3 -- toLiteral-function creates integer literal
4 toLiteral :: Integer -> Expr
5 toLiteral tuid = Lit (ConstI tuid) (TAtom "int4") Nothing
6
7 -- type of provenance sets
8 psetType :: Type

```

## 4.2. PROJECT STRUCTURE

---

```
9  psetType = TArray (TAtom "int4")
10
11  toSet :: [Integer] -> Expr
12  toSet pids = EArrayExpr
13              { typ      = psetType
14                , multidims = False
15                , elements = map toLiteral pids
16                , location  = Nothing
17              }
18
19  -- empty provenance set (identity element)
20  pempty :: Expr
21  pempty = EArrayExpr psetType False [] Nothing
22
23  -- unification of provenance sets
24  (\/) :: Expr -> Expr -> Expr
25  -- static array append using (++)
26  (\/) (EArrayExpr _ _ p1 _) (EArrayExpr _ _ p2 _) =
27      EArrayExpr
28          { typ      = psetType
29            , multidims = False
30            , elements = p1 ++ p2
31            , location  = Nothing
32          }
33  -- dynamic array append (using pg operator)
34  (\/) e1 e2 = EOpexpr { oprleft = Just e1
35                       , oprright = e2
36                       , oprname  = "pg_operator.||"
37                       , typ      = TFunc [] (TPseudo "anyarray")
38                       , location  = Nothing
39                     }
40
41  -- concat for provenance sets
42  pconcat :: [Expr] -> Expr
43  pconcat = foldr (\/) pempty
```

Snippet 4.2: Provenance Sets

## 4.2.6 Push

For each node in the *AST*, we need to be able to add the corresponding provenance annotation to the result. This should work for all expressions, regardless of their type. E.g. in  $\Psi(\mathbf{ROW}(42, \{7\}), \{10\}) = \mathbf{ROW}(42, \{7, 10\})$  the first column is the tuple identifier as expected. Because it does not contain provenance annotations, nothing can be *pushed* into it. However, all other columns contain provenance sets, so the new set  $\{10\}$  is unified with the second column.

In general, *pushing* a label into a data structure means adding it to all parts of the structure where a provenance set can be found. Formally we can define the push-operator  $\Psi$  recursively as follows

$$\Psi(v, \mathbb{P}) := \begin{cases} v \parallel \mathbb{P} & \text{for } v \text{ atomar} \\ \mathbf{ROW}(\rho, \Psi(\text{col}_1, \mathbb{P}), \Psi(\text{col}_2, \mathbb{P}), \dots) & \text{for } v = \mathbf{ROW}(\rho, \text{col}_1, \text{col}_2, \dots) \\ [\Psi(e_1, \mathbb{P}), \Psi(e_2, \mathbb{P}), \dots] & \text{for } v = [e_1, e_2, \dots] \end{cases}$$

If we assume parametric polymorphism, the type of  $\Psi$  looks as follows

$$\Psi : \mathbb{T} \mathbb{P} \rightarrow \mathbb{P} \rightarrow \mathbb{T} \mathbb{P}$$

where  $\mathbb{P}$  is the (static) type of a provenance set and  $\mathbb{T}$  a parameter for an arbitrary *SQL* data structure (e.g. row).

However, since the argument is generally type-polymorphic, *pushing* into it can not be done with just one operator. E.g. the result of a SFW-Block will be a table (with a specific type), while the result of a function application or a column reference might be a literal.

Therefore, the implementation of the  $\Psi$ -operator is particularly challenging. The custom *UDF* definitions and calls for this implementation are presented in the following sections.

### 4.2.6.1 UDF Definitions

Because *SQL* does not support parametric polymorphism though, there can not be one singular definition for  $\Psi$ . Rather, for each specific input type  $\mathbb{T} \mathbb{P}$ , we need to generate a corresponding *UDF*, that also returns  $\mathbb{T} \mathbb{P}$ . These generated *UDF*s

will be written in a log using the **WriterT** monad transformer. Their definitions are dependent on the type. In the scope of this work only atomar values and (unnested) rows are supported. For rows, there needs to be a distinction between custom row types (such as the result of a query or of a **ROW**-constructor) and table values. E.g. for table `flights` the type of  $v$  is `flights` and the corresponding *UDF* definition reflects that (Listing 4.1). The function body reflects the formal definition above: The `tuid` column as primary key is immutable and all other column values are appended using the `||`-operator with the given provenance set.

```
1 CREATE OR REPLACE FUNCTION push(flights2, int4[]) RETURNS
   flights2
2 AS
3 $$
4     SELECT ($1 :: flights2).tuid,
5             (($1 :: flights2).origin) || $2,
6             (($1 :: flights2).destination) || $2,
7             (($1 :: flights2).distance) || $2,
8             (($1 :: flights2).price) || $2
9 $$ LANGUAGE SQL;
```

**Listing 4.1:** UDF definition for table flights

When a type is not already defined in a relation, we need to introduce custom type definitions. This is because, in *UDFs* anonymous composed types are not supported. So, for each composed type, found at a expression with label  $\alpha$ , there is a new type definition whereas  $\alpha$  is added to its name to ensure uniqueness. Considering the example above, if we reduce the type to the `tuid` and `origin` columns, we need to define a new type for the generated *UDF* (Listing 4.2).

```
1 DROP TYPE IF EXISTS row7 CASCADE ;
2 CREATE TYPE row7 AS (tuid int4, origin int4[]);
3 CREATE OR REPLACE FUNCTION push7(row7, int4[]) RETURNS row7
4 AS
5 $$
6     SELECT ($1 :: row7).tuid,
7             (($1 :: row7).origin) || $2
8 $$ LANGUAGE SQL;
```

**Listing 4.2:** UDF and type definition for custom row type



## 4.2.6.2 UDF Calls

As stated above, giving a function definition for  $\Psi$  is not going to be sufficient. Moreover, the function call may have to be wrapped in an additional structure to ensure the output type equals the input type and is valid *SQL*-code. There are two cases, in which function calls need to be wrapped: When pushing into a SFW-block and into a table expression within a **FROM**-clause. A simple cell expression *ec* can just be replaced with the function call. In this generalized example, the arguments for the push operator are flipped for better readability. This format is not valid *SQL* syntax.

```

SELECT var.ρ AS ρ
Ψ3(ℓ3,      Ψ1(ℓ1, ec12) AS col1, ... )
FROM Ψ2(ℓ2, et2) AS var

```

Wrapping the relevant  $\Psi$  function calls each in a subquery, produces this expanded query, which finally is syntactically correct:

```

SELECT call3.*
FROM (
  SELECT var.ρ AS ρ
      Ψ1(ℓ1, ec12) AS col1, ...
  FROM (
    SELECT call2.*
        FROM et2 AS arg2
        Ψ2(ℓ2, arg2) AS call2
    ) AS var
  ) AS arg3
  Ψ3(ℓ3, arg3) AS call3

```

In this example, nothing happens at the call of  $\Psi_1$  because it is a cell expression and the function call returns a cell expression, too. The function call  $\Psi_2$  replacing a table expression is wrapped in a subquery where the function result is dereferenced in the **SELECT**-clause. Finally,  $\Psi_3$  called on a query is wrapped in another query analogous to a table expression.

## 4.3 Translation

The translation rules for Phase1 and Phase2 explained in Section 3.3 are each implemented as a function

$$tr :: \text{Rule } (SQL\ a) \rightarrow (SQL\ a)$$

in the modules Phase1 and Phase2 representing  $\Rightarrow$ . This means that  $tr$  computes for any given  $SQL$  expression a  $SQL$  expression wrapped in a monadic context. In this case the monad is the **OperSem** monad used by the Database Research Group. **OperSem** is a combination of the **ExceptT**, **WriterT**, **StateT** and **ReaderT** monad transformers. Using the `mtl`-package the functionality of this monad stack can be accessed easily without any lifting. In Phase1 the monad stack is only used to throw exceptions, however it is still useful during development to log information about the AST traversal and might have additional purpose in future work. In Phase2 the **WriterT** functionality is used to collect *UDF* and type definitions and for handling exceptions.

Pattern matching the expressions in the AST allows us to project each rule onto a distinct case in the definition of  $tr$ . The function **transformM** exported by the LogParser allows us to 'skip' all nodes which are not relevant for How-Provenance by applying  $tr$  to all subtrees of those nodes. This ensures a full traversal of the AST and makes  $tr$  a total function.

When a translation rule is defined for an expression, the left side of the corresponding rule defines the pattern for the argument of  $tr$ . In this case, using **WriterT** subexpressions are recursively evaluated. For the various types of SFW expressions, plain pattern matching is not powerful enough. Hence, we introduce View Patterns. In this work, these are functions of type **Expr**  $\rightarrow$  **Maybe Expr** that can be applied to pattern expressions in the function definition, where the result is the pattern to be matched with wrapped in a **Just** if the pattern holds the desired properties.

The following uses the rules presented in Section 3.3 to exemplify in further detail how the translation process works. In all of them the language extension **Expr**  $\rightarrow$  **Maybe Expr** is used to bind each record selector inexplicitly to a variable with the construct `{..}`.

### 4.3.1 Operators

The implementation of the operator rule is fairly simple: In `Phase1` subexpressions are recursively evaluated and the operator expression containing them is being returned.

```

1  module Phase1 where
2
3  tr :: Rule (SQL a) (SQL a)
4
5  -- ...
6
7
8  tr E0pexpr {...} = do
9      oprleft <- mapM tr oprleft -- map over Maybe Expr
10     oprright <- tr oprright
11     return $ E0pexpr {...}
12
13  -- ...

```

#### Snippet 4.3: Phase1 Rule: Binary Operator

`Phase2` works similarly, only instead of returning the expression, the provenance set containing both sets of the subexpressions is returned together with the label of the operator expression.

```

1  module Phase2 where
2
3  tr :: Rule (SQL a) (SQL a)
4
5  -- ...
6
7  tr (Lab l E0pexpr {...}) = do
8      oprleft <- mapM tr oprleft -- map over Maybe Expr
9      oprright <- tr oprright
10     --operands are expected to yield an atomic value
11     return $ fromMaybe pempty oprleft \/ oprright \/ pset l
12
13  -- ...

```

#### Snippet 4.4: Phase2 Rule: Binary Operator

### 4.3.2 Case

In this rule, again, subexpressions are evaluated. Additionally, free variables are collected to be included in the logging process.

In Phase1 a new inner **ECASEExpr** is created to form the argument for the logging *UDF*. Therefore, all *when* and *then* expressions are extracted. Finally, new *then* expressions are created to use the integer labels, as shown in Section 3.3.2.

```
1  module Phase1 where
2
3  tr :: Rule (SQL a) (SQL a)
4
5  -- ...
6
7  tr (Lab l ECASEExpr {...}) = do
8    let freeVars = scan ECASEExpr {...}
9    args <- mapM tr args
10   defresult <- tr defresult
11   let whens = map compareExpr args
12   let thens = map result args
13   let innerCase = ECASEExpr
14     { arg = Nothing
15     , args = zipWith createWhenExpr whens [1..]
16     , defresult = Lit (ConstI 0) (TAtom "int4") Nothing
17     , typ = TRow []
18     , location = Nothing
19     }
20   unless (isNothing arg) (throwError "CASE: pattern does not match")
21   let arg = Just $ writeCase $ toLiteral l:innerCase:map tuidRef freeVars
22   let args = zipWith createThenExpr thens [1..]
23   return $ ECASEExpr {...}
24
25  -- ...
```

#### Snippet 4.5: Phase1 Rule: Case

For Phase2 *UDF* and type definitions are created and logged based on the return type of the case expression. The provenance of *when* expressions is culminated and a function to read the log is added to the top of the case expression. Finally, the corresponding label is pushed into the result of the whole expression.

```

1  module Phase2 where
2
3  tr :: Rule (SQL a) (SQL a)
4
5  -- ...
6
7  tr (Lab l ECaseExpr {...}) = do
8      let freeVars = scan ECaseExpr {...}
9      let (udf,call) = push typ l
10     tell $ LogC [udf] $ typeDef typ l
11     args <- mapM tr args
12     defresult' <- tr defresult
13     let whens = map pconcat $ culminate $ map compareExpr args
14     let thens = map result args
15     unless (isNothing arg) (throwError "CASE: pattern does not match.")
16     let arg = Just $ readCase $ toLiteral l:map tuidRef freeVars
17     let args = zipWith3 pushWhenExprs whens thens [1..]
18     let defresult = call [defresult', whyProv (last whens)]
19     return $ call [ECaseExpr {...}, pset l]
20
21  -- ...

```

Snippet 4.6: Phase2 Rule: Case

### 4.3.3 *n*-way-Join

Here, the implementation additionally introduces the usage of the language extension **ViewPatterns** to specify the kind of query block. The pattern only matches with the specifications that need to be ensured by normalization, if this is not the case, no expression will match and the compiler throws an exception.

In both phases again, the expression is scanned for free variables and subexpressions are recursively evaluated.

Phase1 adds the *tuid* column to the *select* clause by passing the label and all relevant tuple identifiers to the logging function. This logging function is then added to the top of the *select* statement, while the rest stays the same.

```

1  {-# LANGUAGE ViewPatterns #-}
2  -- ...
3

```

### 4.3. TRANSLATION

---

```
4 module Phase1 where
5
6 import qualified ViewPatterns as VP
7 -- ...
8
9 tr :: Rule (SQL a) (SQL a)
10
11 -- ...
12
13 -- n-way Join
14 tr (Lab l (VP.nWayJoin -> Just QBlockGeneric {...})) = do
15   let freeVars = scan QBlockGeneric {...}
16     from <- mapM tr from
17     whereEx <- mapM tr whereEx
18     select' <- mapM tr select
19     -- list of all toLiteral columns of all tables involved
20     let colRef = map (tuidRef . rRowname) from
21     let cast = toLiteral l -- convert Haskell label to Literal in postgres
22     -- add writeJoin to select
23     let select = writeJoin (cast:colRef++map tuidRef freeVars):select'
24     return $ QBlockGeneric {...}
25
26 -- ...
```

#### Snippet 4.7: Phase1 Rule: n-way join

In Phase2 the logging function is added to the bottom of the from clause to form a join. The locally bound function *p* is used to add provenance from the *where* clause to all expressions in the *select* clause, while the original *where* clause is nullified. The pushing-process including logging is here abbreviated for query blocks to a singular function *psi*.

```
1 {-# LANGUAGE ViewPatterns #-}
2 -- ...
3
4 module Phase2 where
5
6 import qualified ViewPatterns as VP
7 -- ...
8
9 tr :: Rule (SQL a) (SQL a)
```

```

10
11 -- ...
12
13 -- n-way Join
14 tr (Lab l (VP.nWayJoin -> Just QBlockGeneric {...})) = do
15   let freeVars = scan QBlockGeneric {...}
16   from' <- mapM tr from
17   let colRef = map (tuidRef . rRowname) from' -- get all rownames
18   let cast = toLiteral l -- convert Haskell label to Literal in postgres
19   -- last element of from references readJoin
20   let from = from' ++ [readJoin (cast:colRef++map tuidRef freeVars)]
21   select' <- mapM tr select
22   whereEx' <- mapM tr whereEx
23   -- add atomar Push for each colRef
24   let p ETargetEx {...} = ETargetEx
25     {expr= expr \/ fromMaybe pempty whereEx',...}
26   let select = tuidDeRef:map p select'
27   let whereEx = Nothing
28   psi (QBlockGeneric {...},l) -- push l and log udf/type def.
29
30 -- ...

```

Snippet 4.8: Phase2 Rule: n-way join





## 5 | Outlook & Conclusion

We have successfully implemented the core process required for automating how-provenance. The general goal of this work has been achieved: For the central query constructs of *SQL* a reliable translation can be generated that is executable by a DBMS and returns the provenance data as expected. Additional rules can easily be inserted by simply adding the relevant pattern to the translation function. However, there are various other tasks to be done in possible further work:

*Push* is not yet implemented as a total function because at this moment nested rows can not be dealt with. These require a broader implementation of new defined types for each sub-row and a more case-sensitive recursive definition of the  $\Psi$  operator.

Furthermore, the automatically defined *UDFs* can be subject to a vast optimization by removing duplicate type and function definitions and replacing atomic function definitions with the *SQL*-internal union operator. This can reduce the runtime of Phase2 queries significantly, because *push* operations impose the biggest boilerplate in the generated code.

One prerequisite for the compiler to work reliably is the normalization of input queries. This is particularly important for the distinction of SFW-expressions. Normalizing by hand is not only inconvenient, it also introduces additional how-provenance which is not desired. In further work, by automating the normalization process after annotating the input query, this problem can be evaded.

Furthermore, for a more day-to-day orientated usage, the addition of a graphical user interface is an important step. It may help greatly for debugging because

---

*how*-provenance could be visualized e.g. by hovering over the query result. A lot of editors support similar plugins for different programming languages already, so an extension of this sort can be easily imagined.

Finally, as suggested by the  $\Upsilon$ -operator in the appendix, the translation can be expanded to include *why*-provenance. This way the provenance of input cells, which at this point is empty, would be taken into account, too.

# Bibliography

- [CB74] Donald D. Chamberlin and Raymond F. Boyce. Sequel: A structured english query language. In *Proceedings of the 1974 ACM SIGFIDET (Now SIGMOD) Workshop on Data Description, Access and Control, SIGFIDET '74*, page 249–264, New York, NY, USA, 1974. Association for Computing Machinery. doi:10.1145/800296.811515.
- [CCT09] J. Cheney, L. Chiticariu, and W.C. Tan. *Provenance in Databases: Why, How, and Where*. Foundations and Trends in Databases, Issue 4. Now Publishers, 2009. URL: <https://books.google.de/books?id=0HD42deeWVEC>.
- [GKT07] Todd J. Green, Grigoris Karvounarakis, and Val Tannen. Provenance Semirings. In *Proceedings of the Twenty-Sixth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS '07*, page 31–40, New York, NY, USA, 2007. Association for Computing Machinery. doi:10.1145/1265530.1265535.
- [Hir17] Denis Hirn. Compilation of SQL into KL. B.Sc. Thesis, 2017. URL: <https://db.inf.uni-tuebingen.de/attachments/thesis-hirn-2017.pdf>.
- [Jon03] Simon Peyton Jones. Haskell 98 Languages and Libraries. 2003. URL: <https://www.haskell.org/definition/haskell98-report.pdf>.
- [Lud03] Bertram Ludäscher. A Brief Tour through Provenance in Scientific Workflows and Databases. 2003. URL: <http://hdl.handle.net/2142/89717>.
- [MDG18] Tobias Müller, Benjamin Dietrich, and Torsten Grust. You Say ‘What’, I Hear ‘Where’ and ‘Why’ — (Mis-)Interpreting SQL to Derive Fine-

- Grained Provenance. pages 1536–1549, Rio de Janeiro, Brazil, 2018. URL: <http://www.vldb.org/pvldb/vol11/p1536-muller.pdf>.
- [OMG18] Daniel O’Grady, Tobias Müller, and Torsten Grust. How ‘How’ Explains What ‘What’ Computes — How-Provenance for SQL and Query Compilers. London, UK, 2018. URL: <https://www.usenix.org/system/files/conference/tapp2018/tapp2018-paper-ogrady.pdf>.
- [Par18] Gabriel Paradzik. Language-Level Provenance Analysis of SQL. B.Sc. Thesis, 2018. URL: <https://db.inf.uni-tuebingen.de/attachments/paradzik-2018.pdf>.
- [pos] PostgreSQL Documentation. URL: <https://www.postgresql.org/docs/>.

All *URLs* have been checked on August 13, 2020.

# Appendices



# A | List of Translation Rules

This chapter offers a complete list of all formal translation rules implemented in the scope of this work. As a stepping stone for further work the  $Y$ -operator is included in these rules. With it, data provenance can easily be turned into why-Provenance. In the scope of this work is only how-Provenance, hence for now it is implemented as identity.

## A.1 Literals

$$\begin{array}{c} \text{LITERAL} \\ \ell^\alpha \Rightarrow^{cell} (\ell, \{\alpha\}) \end{array}$$

## A.2 Operators

$$\begin{array}{c} \text{BINOP SCALAR} \\ \frac{ec_1^{\alpha_1} \Rightarrow^{cell} (ec_1^1, ec_1^2) \quad ec_2^{\alpha_2} \Rightarrow^{cell} (ec_2^1, ec_2^2)}{(ec_1^{\alpha_1} \otimes ec_2^{\alpha_2})^\alpha \Rightarrow^{cell} (ec_1^1 \otimes ec_2^1, ec_1^2 \cup ec_2^2 \cup \{\alpha\})} \end{array}$$

## A.3 Column References

$$\begin{array}{c} \text{COLUMN} \\ \frac{ec^\alpha \Rightarrow^{cell} (ec^1, ec^2)}{(ec^\alpha \cdot col^{\alpha_2})^{\alpha_0} \Rightarrow^{cell} (ec_1^1 \cdot col, \Psi(ec_1^2 \cdot col, \{\alpha_0, \alpha_2\}))} \end{array}$$

## A.4 Case

$$\begin{array}{c}
\text{CASE} \\
\left| ec_{w_i}^{\alpha_{w_i}} \Rightarrow^{cell} (ec_{w_i}^1, ec_{w_i}^2) \right|_{i=1\dots} \quad \left| ec_{t_i}^{\alpha_{t_i}} \Rightarrow^{cell} (ec_{t_i}^1, ec_{t_i}^2) \right|_{i=0\dots} \\
\{f_1, \dots\} := \text{fv}(e_{in}) \quad \mathcal{O} := \text{tolocation}(\alpha) \\
\text{CASE writeCase}(\mathcal{O}, \\
\quad \text{CASE} \\
\quad \quad \text{WHEN } ec_{w_1}^1 \text{ THEN } 1 \\
\quad \quad \vdots \\
\quad \quad \text{ELSE } 0 \\
X^1 := \quad \text{END} \\
\quad , f_1 \cdot \rho, \dots) \\
\quad \text{WHEN } 1 \text{ THEN } ec_{t_1}^1 \\
\quad \quad \vdots \\
\quad \quad \text{ELSE } ec_{t_0}^1 \\
\quad \text{END} \\
\text{CASE readCase}(\mathcal{O}, f_1 \cdot \rho, \dots) \\
\quad \text{WHEN } 1 \text{ THEN } \Psi(ec_{t_1}^2, Y(ec_{w_1}^2)) \\
\quad \quad \vdots \\
X^2 := \quad \text{WHEN } i \text{ THEN } \Psi(ec_{t_i}^2, Y(ec_{w_1}^2 \cup \dots \cup ec_{w_i}^2)) \\
\quad \quad \vdots \\
\quad \quad \text{ELSE } \Psi(ec_{t_0}^2, Y(ec_{w_1}^2 \cup \dots)) \\
\quad \text{END} \\
\hline
e_{in} := \left( \begin{array}{c} \text{CASE} \\ \quad \text{WHEN } ec_{w_1}^{\alpha_{w_1}} \text{ THEN } ec_{t_1}^{\alpha_{t_1}} \\ \quad \quad \vdots \\ \quad \quad \text{ELSE } ec_{t_0}^{\alpha_{t_0}} \\ \text{END} \end{array} \right)^\alpha \Rightarrow^{cell} (X^1, \Psi(X^2, \{\alpha\}))
\end{array}$$

## A.5 Row Constructors

$$\begin{array}{c}
\text{Row} \\
\left| ec_i^{\alpha_i} \Rightarrow^{cell} (ec_i^1, ec_i^2) \right|_{i=1\dots} \quad \rho := \text{totuid}(\alpha_0) \\
\hline
(\text{ROW}(ec_1^{\alpha_1}, \dots))^\alpha \Rightarrow^{cell} (\text{ROW}(\rho, ec_1^1, \dots), \Psi(\text{ROW}(\rho, ec_1^2, \dots), \{\alpha\}))
\end{array}$$



## A.6 Values

VALUES

$$\frac{\left| ec_i^{\alpha_i} \Rightarrow^{cell} (ec_i^1, ec_i^2) \right|_{i=1\dots}}{(\mathbf{VALUES} (ec_1^{\alpha_1}, \dots))^{\alpha} \Rightarrow^{table} (\mathbf{VALUES} (ec_1^1, \dots), \Psi(\mathbf{VALUES} (ec_1^2, \dots), \{\alpha\}))}$$

## A.7 Union All

UNION ALL

$$\begin{aligned} & \mathbb{O} := \text{tolocation}(\alpha) \quad \{f_1, \dots\} := \text{fv}(et_1^{\alpha_1}) \\ & \quad [col_1, \dots] := \text{columns}(et_1) \\ & et_1^{\alpha_1} \Rightarrow^{table} (et_1^1, et_1^2) \quad et_2^{\alpha_2} \Rightarrow^{table} (et_2^1, et_2^2) \\ X^1 := & \quad \mathbf{SELECT} \text{ writeUnion}(\mathbb{O}, var.\rho, f_1.\rho, , \dots) \mathbf{AS} \rho, \\ & \quad var.col_1, \dots \\ & \quad \mathbf{FROM} et_1^1 \mathbf{AS} var \\ X^2 := & \quad \mathbf{SELECT} \text{ readUnion}(\mathbb{O}, var.\rho, f_1.\rho, , \dots) \mathbf{AS} \rho, \\ & \quad var.col_1, \dots \\ & \quad \mathbf{FROM} et_1^2 \mathbf{AS} var \end{aligned}$$

$$\frac{}{(et_1^{\alpha_1} \mathbf{UNION ALL} et_2^{\alpha_2})^{\alpha} \Rightarrow^{table} (X^1 \mathbf{UNION ALL} et_2^1, \Psi(X^2 \mathbf{UNION ALL} et_2^2, \{\alpha\}))}$$

## A.8 With

WITH

$$\frac{\left| et_i^{\alpha_i} \Rightarrow^{table} (et_i^1, et_i^2) \right|_{i=0,\dots}}{X^1 := \mathbf{WITH} (tab1_1 \mathbf{AS} et_1^1, \dots) et_0^1 \\ X^2 := \mathbf{WITH} (tab2_1 \mathbf{AS} et_1^2, \dots) et_0^2} \\ (\mathbf{WITH} (tab_1 \mathbf{AS} et_1^{\alpha_1}, \dots) et_0^{\alpha_0})^{\alpha} \Rightarrow^{table} (X^1, \Psi(X^2, \{\alpha\}))$$

## A.9 SFW: Plain Select From

SFW-MAP

$$\begin{array}{l}
\left| \begin{array}{l}
ec_i^{\alpha_i} \Rightarrow^{cell} (ec_i^1, ec_i^2) \\
et^{\alpha_0} \Rightarrow^{table} (et^1, et^2)
\end{array} \right|_{i=1..n} \\
X^1 := \text{SELECT } var.\rho \text{ AS } \rho \\
\qquad \qquad \qquad ec_1^1 \text{ AS } col_1, \dots \\
\qquad \qquad \qquad \text{FROM } et^1 \text{ AS } var \\
X^2 := \text{SELECT } var.\rho \text{ AS } \rho \\
\qquad \qquad \qquad ec_1^2 \text{ AS } col_1, \dots \\
\qquad \qquad \qquad \text{FROM } et^2 \text{ AS } var \\
\hline
\left( \begin{array}{l}
\text{SELECT } ec_1^{\alpha_1} \text{ AS } col_1, \dots \\
\text{FROM } et^{\alpha_0} \text{ AS } var
\end{array} \right)^\alpha \Rightarrow^{table} (X^1, \Psi(X^2, \{\alpha\}))
\end{array}$$

## A.10 SFW: $n$ -way join

SFW-JOIN

$$\begin{array}{l}
\mathbb{O} := \text{cast}(\alpha) \quad \{f_1, \dots\} := \text{fv}(et_{inp}) \\
\left| \begin{array}{l}
ec_i^{\alpha_i} \Rightarrow^{cell} (ec_i^1, ec_i^2) \\
et_i^{\alpha_i} \Rightarrow^{table} (et_i^1, et_i^2)
\end{array} \right|_{i=0, \dots, n} \\
\qquad \qquad \qquad \left| \qquad \qquad \qquad \right|_{i=(n+1), \dots, (n+m)} \\
X^1 := \text{SELECT } \text{writeJoin}(\mathbb{O}, var_1.\rho, \dots, f_1.\rho, \dots) \text{ AS } \rho \\
\qquad \qquad \qquad ec_1^1 \text{ AS } col_1, \dots \\
\qquad \qquad \qquad \text{FROM } et_1^1 \text{ AS } var_1, \dots \\
\qquad \qquad \qquad \text{WHERE } ec_0^1 \\
X^2 := \text{SELECT } log.\rho \text{ AS } \rho \\
\qquad \qquad \qquad \Psi(ec_1^2, Y(ec_0^2)) \text{ AS } col_1, \dots \\
\qquad \qquad \qquad \text{FROM } et_1^2 \text{ AS } var_1, \dots, \\
\qquad \qquad \qquad \text{readJoin}(\mathbb{O}, var_1.\rho, \dots, f_1.\rho, \dots) \text{ AS } log \\
\hline
et_{inp} := \left( \begin{array}{l}
\text{SELECT } ec_1^{\alpha_1} \text{ AS } col_1, \dots \\
\text{FROM } et_{n+1}^{\alpha_{n+1}} \text{ AS } var_{n+1}, \dots \\
\text{WHERE } ec_0^{\alpha_0}
\end{array} \right)^\alpha \Rightarrow^{table} (X^1, \Psi(X^2, \{\alpha\}))
\end{array}$$

## A.11 SFW: Aggregates

SFW-AGG

$$\begin{array}{l}
 \textcircled{\alpha} := \text{tolocation}(\alpha) \quad et^{\alpha_0} \Rightarrow^{table} (et^1, et^2) \quad \{f_1, \dots\} := \text{fv}(et_{inp}) \\
 \left| ec_i^{\alpha_i} \Rightarrow^{cell} (ec_i^1, ec_i^2) \right|_{i=1..n} \quad \left| ec_{gr_i}^{\alpha_{gr_i}} \Rightarrow^{cell} (ec_{gr_i}^1, ec_{gr_i}^2) \right|_{i=1..n} \\
 \\
 \text{SELECT writeAgg}(\textcircled{\alpha}, \text{ARRAY\_AGG}(var.\rho), f_1.\rho, \dots) \text{ AS } \rho, \\
 \\
 X^1 := \quad \oplus(ec_1^1) \text{ AS } col_1, \dots \\
 \text{FROM } et^1 \text{ AS } var \\
 \text{GROUP BY } ec_{gr_1}^1, \dots \\
 \\
 \text{SELECT THE}(log.\rho) \text{ AS } \rho, \\
 \quad \bigcup \Psi(ec_1^2, Y(ec_{gr_1}^2 \cup \dots)) \text{ AS } col_1, \dots \\
 X^2 := \quad \text{FROM } et^2 \text{ AS } var, \\
 \quad \text{readAgg}(\textcircled{\alpha}, var.\rho, f_1.\rho, \dots) \text{ AS } log \\
 \quad \text{GROUP BY } log.\rho \\
 \hline
 et_{inp} := \left( \begin{array}{l} \text{SELECT } \oplus(ec_1^{\alpha_1})^{\alpha_{agg1}} \text{ AS } col_1, \dots \\ \text{FROM } et^{\alpha_0} \text{ AS } var \\ \text{GROUP BY } ec_{gr_1}^{\alpha_{gr_1}}, \dots \end{array} \right)^\alpha \Rightarrow^{table} (X^1, \Psi(X^2, \{\alpha\}))
 \end{array}$$

## A.12 SFW: Order By

SFW-ORDERBY

$$\mathbb{O} := \text{tolocation}(\alpha) \quad et^{\alpha_0} \Rightarrow^{table} (et^1, et^2) \quad \{f_1, \dots\} := \text{fv}(et_{inp})$$

$$\left| ec_i^{\alpha_i} \Rightarrow^{cell} (ec_i^1, ec_i^2) \right|_{i=1..} \quad \left| ec_{dis_i}^{\alpha_{dis_i}} \Rightarrow^{cell} (ec_{dis_i}^1, ec_{dis_i}^2) \right|_{i=1..}$$

**SELECT** writeFilter ( $\mathbb{O}, var''.\rho, f_1.\rho, \dots$ ) **AS**  $\rho$   
 $var''.$ col<sub>1</sub> **AS** col<sub>1</sub>, ...

**FROM** (

**SELECT**  $var'$ . $\rho$  **AS**  $\rho$   
 $ec_1^1$  **AS** col<sub>1</sub>, ...

**FROM**  $et^1$  **AS**  $var'$

**ORDER BY** col<sub>ord<sub>1</sub></sub> **ASC | DESC**, ...

**DISTINCT ON**  $ec_{dis_1}^1$ , ...

**LIMIT**  $\ell_l$

**OFFSET**  $\ell_o$

) **AS**  $var''$

**ORDER BY**  $var''.$ col<sub>ord<sub>1</sub></sub> **ASC | DESC**, ...

$X^1 :=$

**SELECT** log. $\rho$  **AS**  $\rho$

$\Psi(ec_1^2, p)$  **AS** col<sub>1</sub>, ...

$X^2 :=$

**FROM**  $et^2$  **AS**  $var, \dots,$

readFilter ( $\mathbb{O}, var.\rho, f_1.\rho, \dots$ ) **AS** log

$$p := \{\alpha\} \cup Y(\{\alpha_{ord_1}\} \cup \dots \cup var.col_{ord_1} \cup \dots \cup ec_{dis_1}^2 \cup \dots)$$

$$et_{inp} := \left( \begin{array}{l} \mathbf{SELECT} \ ec_1^{\alpha_1} \ \mathbf{AS} \ col_1, \dots \\ \mathbf{FROM} \ et^{\alpha_0} \ \mathbf{AS} \ var \\ \mathbf{ORDER \ BY} \ col_{ord_1}^{\alpha_{ord_1}} \ \mathbf{ASC} \ | \ \mathbf{DESC} \ , \dots \\ \mathbf{DISTINCT \ ON} \ ec_{dis_1}^{\alpha_{dis_1}} \ , \dots \\ \mathbf{LIMIT} \ \ell_l \\ \mathbf{OFFSET} \ \ell_o \end{array} \right)^\alpha \Rightarrow^{table} (X^1, X^2)$$

## B | Sample Output

The following listing shows the output of the compiler presented in this work with the query in Listing 1.1 as input. This generated *SQL* code can be executed by Postgres and results in the desired Provenance data.

```
1  -- view definitions
2  CREATE MATERIALIZED VIEW flights1 AS
3      SELECT nextval('tuid_seq') AS tuid,
4             flights.origin AS origin,
5             flights.destination AS destination,
6             flights.distance AS distance,
7             flights.price AS price
8      FROM flights;
9  CREATE MATERIALIZED VIEW flights2 AS
10     SELECT flights1.tuid AS tuid,
11            ARRAY[] :: int4[] AS origin,
12            ARRAY[] :: int4[] AS destination,
13            ARRAY[] :: int4[] AS distance,
14            ARRAY[] :: int4[] AS price
15     FROM flights1;
16
17  -- phase 1 - query:
18  (SELECT writeunion(20, subquery1.tuid :: int4) AS tuid,
19     subquery1.destination AS destination
20   FROM (SELECT writejoin(8, RTE0.tuid :: int4) AS tuid,
21          RTE0.destination AS destination
22         FROM flights1 AS RTE0(tuid,
23                                origin,
24                                destination,
25                                distance,
26                                price)
27        WHERE RTE0.price < 100
28       ) AS subquery1(tuid, destination))
29  UNION ALL
30  (SELECT writejoin(17, RTE2.tuid :: int4) AS tuid,
31     RTE2.destination AS destination
32   FROM flights1 AS RTE2(tuid,
33                          origin,
34                          destination,
```

```

35         distance,
36         price)
37 WHERE RTE2.price > 200);
38
39 -- phase 2 - udf and type definitions:
40 DROP TYPE IF EXISTS row4 CASCADE;
41 CREATE TYPE row4 AS (tuid int4, origin int4[], destination int4[], distance int4[],
42     price int4[]);
43 DROP TYPE IF EXISTS row8 CASCADE;
44 CREATE TYPE row8 AS (tuid int4, destination int4[]);
45 DROP TYPE IF EXISTS row10 CASCADE;
46 CREATE TYPE row10 AS (tuid int4, destination int4[]);
47 DROP TYPE IF EXISTS row13 CASCADE;
48 CREATE TYPE row13 AS (tuid int4, origin int4[], destination int4[], distance int4[],
49     price int4[]);
50 DROP TYPE IF EXISTS row17 CASCADE;
51 CREATE TYPE row17 AS (tuid int4, destination int4[]);
52 DROP TYPE IF EXISTS row19 CASCADE;
53 CREATE TYPE row19 AS (tuid int4, destination int4[]);
54 DROP TYPE IF EXISTS row21 CASCADE;
55 CREATE TYPE row21 AS (tuid int4, destination int4[]);
56 CREATE OR REPLACE FUNCTION push4(flights2, int4[]) RETURNS flights2
57 AS
58 $$
59     SELECT ($1 :: flights2).tuid,
60            (($1 :: flights2).origin) || $2,
61            (($1 :: flights2).destination) || $2,
62            (($1 :: flights2).distance) || $2,
63            (($1 :: flights2).price) || $2
64 $$ LANGUAGE SQL;
65 CREATE OR REPLACE FUNCTION push2(int4[], int4[]) RETURNS int4[]
66 AS
67 $$
68     SELECT $1 || $2
69 $$ LANGUAGE SQL;
70 CREATE OR REPLACE FUNCTION push5(int4[], int4[]) RETURNS int4[]
71 AS
72 $$
73     SELECT $1 || $2
74 $$ LANGUAGE SQL;
75 CREATE OR REPLACE FUNCTION push8(row8, int4[]) RETURNS row8
76 AS
77 $$
78     SELECT ($1 :: row8).tuid, (($1 :: row8).destination) || $2
79 $$ LANGUAGE SQL;
80 CREATE OR REPLACE FUNCTION push10(row10, int4[]) RETURNS row10
81 AS
82 $$
83     SELECT ($1 :: row10).tuid, (($1 :: row10).destination) || $2
84 $$ LANGUAGE SQL;
85 CREATE OR REPLACE FUNCTION push13(flights2, int4[]) RETURNS flights2

```

```

84 AS
85 $$
86     SELECT ($1 :: flights2).tuid,
87            (($1 :: flights2).origin) || $2,
88            (($1 :: flights2).destination) || $2,
89            (($1 :: flights2).distance) || $2,
90            (($1 :: flights2).price) || $2
91 $$ LANGUAGE SQL;
92 CREATE OR REPLACE FUNCTION push11(int4[], int4[]) RETURNS int4[]
93 AS
94 $$
95     SELECT $1 || $2
96 $$ LANGUAGE SQL;
97 CREATE OR REPLACE FUNCTION push14(int4[], int4[]) RETURNS int4[]
98 AS
99 $$
100    SELECT $1 || $2
101 $$ LANGUAGE SQL;
102 CREATE OR REPLACE FUNCTION push17(row17, int4[]) RETURNS row17
103 AS
104 $$
105    SELECT ($1 :: row17).tuid, (($1 :: row17).destination) || $2
106 $$ LANGUAGE SQL;
107 CREATE OR REPLACE FUNCTION push19(row19, int4[]) RETURNS row19
108 AS
109 $$
110    SELECT ($1 :: row19).tuid, (($1 :: row19).destination) || $2
111 $$ LANGUAGE SQL;
112 CREATE OR REPLACE FUNCTION push21(row21, int4[]) RETURNS row21
113 AS
114 $$
115    SELECT ($1 :: row21).tuid, ($1 :: row21).destination || $2
116 $$ LANGUAGE SQL;
117
118 -- phase 2 - query:
119 SELECT call21.*
120 FROM ((SELECT readunion(20, subquery1.tuid :: int4) AS tuid,
121        subquery1.destination AS destination
122        FROM (SELECT call10.*
123              FROM (SELECT call8.*
124                    FROM (SELECT log.tuid AS tuid,
125                          (push2(RTE0.destination, ARRAY[2] :: int4[]))
126                          ||
127                          (((push5(RTE0.price, ARRAY[5] :: int4[]))
128                          ||
129                          (ARRAY[6] :: int4[]))
130                          ||
131                          (ARRAY[7] :: int4[])) AS destination
132                    FROM (SELECT call4.*
133                          FROM flights2 AS arg4(tuid,
134

```

```

135         destination,
136         distance,
137         price),
138         LATERAL push4(arg4,
139         ARRAY[4] :: int4[]) AS call4(tuid,
140         origin
141     ,
142     destination,
143     distance,
144     price)
145     ) AS RTE0,
146     readjoin(8, RTE0.tuid :: int4) AS log
147     ) AS arg8,
148     LATERAL push8(arg8, ARRAY[8] :: int4[]) AS call8(tuid,
149     destination)
150     ) AS arg10(tuid, destination),
151     LATERAL push10(arg10, ARRAY[10] :: int4[]) AS call10(tuid,
152     destination)
153     ) AS subquery1)
154     UNION ALL
155     (SELECT call19.*
156     FROM (SELECT call17.*
157     FROM (SELECT log.tuid AS tuid,
158     (push11(RTE2.destination, ARRAY[11] :: int4[]))
159     ||
160     ((push14(RTE2.price, ARRAY[14] :: int4[]))
161     ||
162     (ARRAY[15] :: int4[]))
163     ||
164     (ARRAY[16] :: int4[])) AS destination
165     FROM (SELECT call13.*
166     FROM flights2 AS arg13(tuid,
167     origin,
168     destination,
169     distance,
170     price),
171     LATERAL push13(arg13,
172     ARRAY[13] :: int4[]) AS call13(tuid,
173     origin,
174     destination,
175     distance,
176     price)
177     ) AS RTE2,
178     readjoin(17, RTE2.tuid :: int4) AS log
179     ) AS arg17,
180     LATERAL push17(arg17, ARRAY[17] :: int4[]) AS call17(tuid,
181     destination)

```



```
181     ) AS arg19(tuid, destination),
182     LATERAL push19(arg19, ARRAY[19] :: int4[]) AS call19(tuid,
183                                     destination))
184 ) AS arg21,
185 LATERAL push21(arg21,
186               ARRAY[21] :: int4[]) AS call21(tuid,destination);
```

Listing B.1: Output of the compiler for **UNION ALL** query



## Acknowledgements

Foremost, I would like to express my gratitude to my supervisor Dr. Tobias Müller for his continuous support and for providing and explaining the theoretical input for this work, most notably the translation rules.

Moreover I have to thank Denis Hirn for the implementation of the *LogParser* and his support for understanding the usage of his libraries and some general Haskell-related topics.

Lastly, my gratitude goes to my brother Florian Engel and my friend Holly Ransom for proof reading this work.



## Author's Declaration

Hiermit erkläre ich, dass ich die Arbeit selbständig verfasst habe, keine anderen als die angegebenen Hilfsmittel und Quellen benutzt habe, alle wörtlich oder sinngemäß aus anderen Werken übernommenen Aussagen als solche gekennzeichnet habe und dass die Arbeit weder vollständig noch in wesentlichen Teilen Gegenstand eines anderen Prüfungsverfahrens gewesen ist und dass ich die Arbeit weder vollständig noch in wesentlichen Teilen bereits veröffentlicht habe sowie dass das in Dateiform eingereichte Exemplar mit eingereichten gebundenen Exemplaren übereinstimmt.

---

Datum, Ort, Unterschrift