

# A Note on the Translation of SQL to Tuple Calculus

martin gogolla

TU Braunschweig, Informatik, Abt. Datenbanken  
Postfach 3329, D-3300 Braunschweig  
e-mail: gogolla at infbs.uucp

**Abstract:** This note presents a translation of a subset of the relational query language SQL into the well known tuple calculus. Roughly speaking, tuple calculus corresponds to first order predicate calculus. The SQL subset is relationally complete and represents a "relational core" of the language. Nevertheless, our translation is simple and elegant. Therefore it is especially well suited as a beginners course into the principles of a formal definition of SQL.

**Key words:** Database theory, relational query language, tuple calculus, SQL, formal semantics.

## 1. Introduction

Several authors have studied the formal semantics of the relational query language SQL. In [CG 85] SQL is translated into relational algebra, [Bü 87] takes special care of aggregate functions and in [PBG 89] a translation of SQL into tuple calculus is proposed. A comprehensive description and constructive criticism of the SQL standard can be found in [Da 87]. Nevertheless, up to now no elegant proposal has appeared which especially allows to prove properties of the SQL language easily. This is possible within our approach.

But before discussing the details of the translation we point out the need of formal semantics by means of an example. Consider the following simple cooking book database (keys are underlined).

```
RECIPE ( RName : string;  
         Inged : string;  
         UsedQuantity : real )  
STORE ( Inged : string;  
        StoreQuantity : real )  
ORIGIN ( RName : string;  
         Country : string )
```

Let us now look at a typical SQL query looking for some recipes using garlic.

```
SELECT r.RName  
FROM RECIPE r  
WHERE r.Inged='garlic' AND  
r.UsedQuantity >  
ANY ( SELECT r'.UsedQuantity  
      FROM RECIPE r',  
      ORIGIN o  
      WHERE r'.RName=o.RName AND  
            o.Country='italy' AND  
            r'.Inged='garlic' )
```

The subquery computes all garlic quantities in Italian recipes. Therefore, the SQL keywords suggest that the query performs the following task: "Select recipes which use more garlic than any Italian recipe."

In usual colloquial English it is understood that these are recipes using more garlic than the quantity of garlic used in every Italian recipe, i.e., more than the maximal quantity. But a look at the formal semantics shows that the SQL query does not do this job:

```

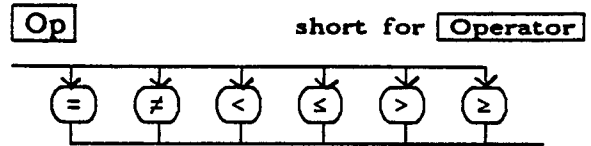
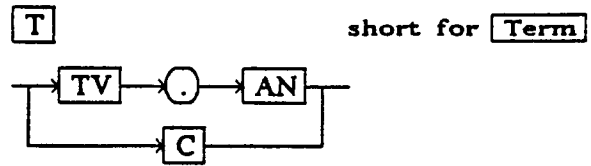
{ res : (Res1) |
∃ r:RECIPE
  res.Res1=r.RName ∧
  r.Ingred='gralic' ∧
  ∃ r':RECIPE
    (∃ o:ORIGIN
      r'.RName=o.RName ∧
      o.Country='italy' ∧
      r'.Ingred='garlic') -
    ∧ r.UsedQuantity>r'.UsedQuantity }

```

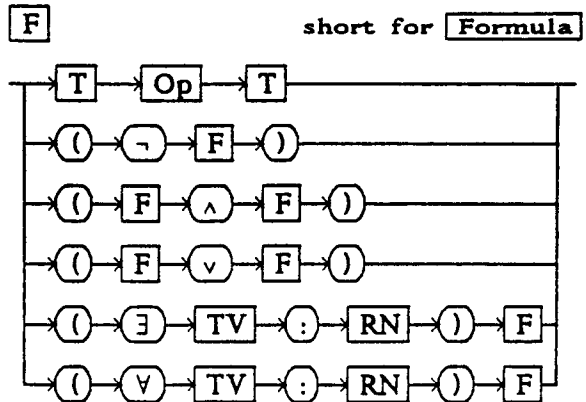
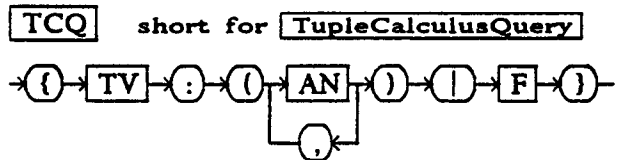
The formula now points out that the SQL query finds recipes which use more garlic than at least one Italian recipe or in other words recipes which use more garlic than the minimal quantity of garlic in Italian recipes. In fact, a similar query to the above one together with the "wrong" description in colloquial English can be found in IBM manuals for SQL [Da 87]. Thus, it is desirable to formally define the semantics of SQL as a basis for implementations. On the other hand, a formal definition can also be employed to prove the equivalence of queries thus making formally well founded optimization techniques even more attractive. Furthermore, one can point out properties of the SQL language.

## 2. Tuple calculus and SQL

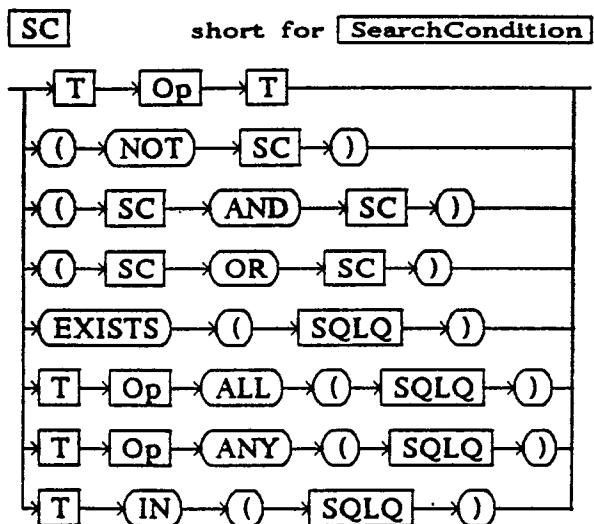
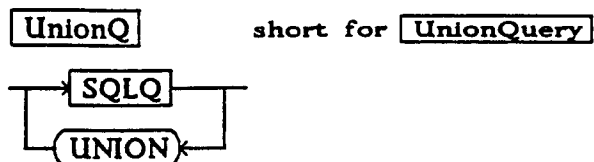
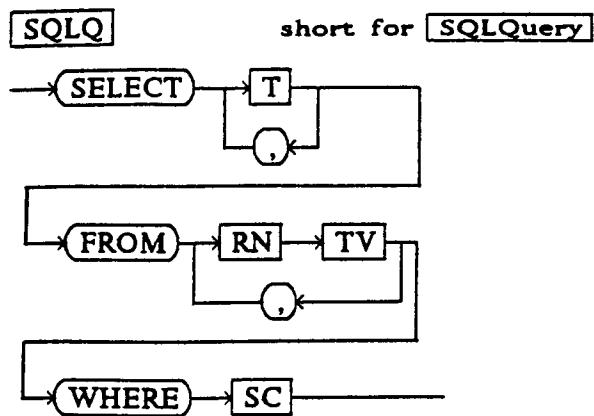
We describe the syntax of our versions of tuple calculus and SQL with syntax diagrams. The semantics of tuple calculus can be found in [Ma 83]; the semantics of SQL is subject to this note and studied in chapter 3. Both, tuple calculus and SQL involve terms  $\bar{T}$  (consisting of tuple variables  $\bar{TV}$  and attribute names  $\bar{AN}$  or constants  $\bar{C}$ ) and comparison operators  $\bar{Op}$ .



As mentioned before, our version of tuple calculus corresponds to first order predicate calculus. The qualifying formula  $\bar{F}$  of the tuple calculus query  $\bar{TCQ}$  has exactly one free variable, namely the specified tuple variable  $\bar{TV}$  whose structure is determined by various attribute names  $\bar{AN}$ .  $\bar{RN}$  stands for the name of a relation.



Our SQL queries  $\bar{SQLQ}$  take into account that subqueries can be built via the operators EXISTS, ALL, ANY and IN specified in the search condition  $\bar{SC}$ . The operator UNION is allowed only on the top level.  $\bar{UnionQ}$  generates these queries.



The nonterminals **TV**, **AN**, **RN** and **C** are left unspecified. Additionally, parenthesis may be omitted in accordance with the usual rules and the tuple calculus also employs implication  $\Rightarrow$  with the usual semantics.

### 3. The translation

There are several general requirements the SQL queries have to fulfill ( $\tau$ ,  $\tau_1$  and  $\tau_2$  stand for terms and  $\varphi$ ,  $\varphi_1$  and  $\varphi_2$  for formulas):

- We assume that the given SQL query "SELECT  $\tau_1, \dots, \tau_n$  FROM  $R_1 s_1, \dots, R_m s_m$  WHERE  $\varphi$ " and all subqueries are qualified completely via explicitly declared tuple variables: If there are references to an attribute  $A_j$  in a result term  $\tau_i$  or in the formula  $\varphi$ , they are of the form  $s_1.A_j$ , where  $s_1$  is declared in the FROM clause (or in the case of a subquery perhaps in the FROM part of an outer query) and  $A_j$  is an attribute of the corresponding relation schema  $R_1$ . The names of the tuple variables must be unique ( $i \neq j \Rightarrow s_i \neq s_j$ ) and must be declared only once. The result terms as well as the formula  $\varphi$  use only declared tuple variables.
- Additionally, the result terms  $\tau_i$  may be constants  $c_i$ . The result variable of the tuple calculus expression must be different from all declared variables in the query.
- If terms are compared with an operator like " $\tau_1 \omega \tau_2$ " or " $\tau_1 \omega \text{ALL} (\text{SELECT } \tau_2 \text{ FROM } \dots \text{ WHERE } \dots)$ ", then  $\tau_1$  and  $\tau_2$  have the same data type. For a UNION expression we assume  $\tau_1$  and  $\tau_i'$  have the same data type for  $i \in 1..n$ . The translation is now done by defining a function *sql2tc*, which gives for every SQL expression the corresponding tuple calculus one.

$$\begin{aligned}
 \text{sql2tc} \llbracket \text{SELECT } \tau_1, \dots, \tau_n \\
 \text{FROM } R_1 s_1, \dots, R_m s_m \\
 \text{WHERE } \varphi \rrbracket := \\
 \{ r : ( \text{Res}_1, \dots, \text{Res}_n ) \mid \\
 ( \exists s_1:R_1, \dots, s_m:R_m ) \\
 ( r.\text{Res}_1 = \tau_1 \wedge \dots \wedge r.\text{Res}_n = \tau_n \\
 \wedge \text{sql2tc} \llbracket \varphi \rrbracket ) \} \\
 \\
 \text{sql2tc} \llbracket ( \text{NOT } \varphi ) \rrbracket := \\
 ( \neg \text{sql2tc} \llbracket \varphi \rrbracket )
 \end{aligned}$$

$sql2tc \llbracket ( \varphi_1 \text{ AND } \varphi_2 ) \rrbracket :=$   
 $( sql2tc \llbracket \varphi_1 \rrbracket \wedge sql2tc \llbracket \varphi_2 \rrbracket )$   
 $sql2tc \llbracket ( \varphi_1 \text{ OR } \varphi_2 ) \rrbracket :=$   
 $( sql2tc \llbracket \varphi_1 \rrbracket \vee sql2tc \llbracket \varphi_2 \rrbracket )$   
 $sql2tc \llbracket \tau_1 \omega \tau_2 \rrbracket := \tau_1 \omega \tau_2$   
 $sql2tc \llbracket \tau \omega \text{ ALL}$   
 $( \text{ SELECT } s_1.A$   
 $\text{ FROM } R_1 s_1, \dots, R_m s_m$   
 $\text{ WHERE } \varphi ) \rrbracket :=$   
 $(\forall s_1:R_1)$   
 $((\exists s_1:R_1, \dots, s_{i-1}:R_{i-1}, s_{i+1}:R_{i+1}, \dots, s_m:R_m)$   
 $sql2tc \llbracket \varphi \rrbracket ) \Rightarrow \tau \omega s_1.A )$   
 $sql2tc \llbracket \tau \omega \text{ ANY}$   
 $( \text{ SELECT } s_1.A$   
 $\text{ FROM } R_1 s_1, \dots, R_m s_m$   
 $\text{ WHERE } \varphi ) \rrbracket :=$   
 $(\exists s_1:R_1)$   
 $((\exists s_1:R_1, \dots, s_{i-1}:R_{i-1}, s_{i+1}:R_{i+1}, \dots, s_m:R_m)$   
 $sql2tc \llbracket \varphi \rrbracket ) \wedge \tau \omega s_1.A )$   
 $sql2tc \llbracket \tau \text{ IN}$   
 $( \text{ SELECT } s_1.A$   
 $\text{ FROM } R_1 s_1, \dots, R_m s_m$   
 $\text{ WHERE } \varphi ) \rrbracket :=$   
 $(\exists s_1:R_1)$   
 $((\exists s_1:R_1, \dots, s_{i-1}:R_{i-1}, s_{i+1}:R_{i+1}, \dots, s_m:R_m)$   
 $sql2tc \llbracket \varphi \rrbracket ) \wedge \tau = s_1.A )$   
 $sql2tc \llbracket \text{ EXISTS}$   
 $( \text{ SELECT } r_1.A_1, \dots, r_n.A_n$   
 $\text{ FROM } R_1 s_1, \dots, R_m s_m$   
 $\text{ WHERE } \varphi ) \rrbracket :=$   
 $(\exists s_1:R_1, \dots, s_m:R_m) sql2tc \llbracket \varphi \rrbracket$   
 $sql2tc \llbracket \text{ SELECT } \tau_1, \dots, \tau_n$   
 $\text{ FROM } R_1 s_1, \dots, R_m s_m$   
 $\text{ WHERE } \varphi$   
 $\text{ UNION}$   
 $\text{ SELECT } \tau_1', \dots, \tau_n'$   
 $\text{ FROM } R_1' s_1', \dots, R_k' s_k'$   
 $\text{ WHERE } \varphi' \rrbracket :=$

$\{ r : ( \text{ Res}_1, \dots, \text{ Res}_n ) \mid$   
 $( \exists s_1:R_1, \dots, s_m:R_m )$   
 $( r.\text{Res}_1 = \tau_1 \wedge \dots \wedge r.\text{Res}_n = \tau_n$   
 $\wedge sql2tc \llbracket \varphi \rrbracket )$   
 $\vee$   
 $( \exists s_1':R_1', \dots, s_k':R_k' )$   
 $( r.\text{Res}_1 = \tau_1' \wedge \dots \wedge r.\text{Res}_n = \tau_n'$   
 $\wedge sql2tc \llbracket \varphi' \rrbracket ) \}$

The rule for UNION can be generalized to the case with more than 2 operands. All used tuple variables are of the form  $r:R$ ; therefore we only employ safe expressions [Ma 83].

#### 4. Some Laws

On the basis of a formally defined language it is now possible to state properties of the language. Here are some laws which can be proved formally.

##### Theorem:

- (I)  $\tau_1 \text{ IN } ( \text{ SELECT } \tau_2$   
 $\text{ FROM } \rho$   
 $\text{ WHERE } \varphi ) \Leftrightarrow$   
 $\tau_1 = \text{ ANY } ( \text{ SELECT } \tau_2$   
 $\text{ FROM } \rho$   
 $\text{ WHERE } \varphi )$
- (II)  $\text{ NOT } ( \tau_1 \omega \text{ ALL } ( \text{ SELECT } \tau_2$   
 $\text{ FROM } \rho$   
 $\text{ WHERE } \varphi ) ) \Leftrightarrow$   
 $\tau_1 \bar{\omega} \text{ ANY } ( \text{ SELECT } \tau_2$   
 $\text{ FROM } \rho$   
 $\text{ WHERE } \varphi )$
- (III)  $\tau_1 \omega \text{ ANY } ( \text{ SELECT } \tau_2$   
 $\text{ FROM } \rho$   
 $\text{ WHERE } \varphi ) \Leftrightarrow$   
 $\text{ EXISTS } ( \text{ SELECT } \tau(\rho)$   
 $\text{ FROM } \rho$   
 $\text{ WHERE } (\varphi) \text{ AND } \tau_1 \omega \tau_2 )$
- (IV)  $\tau_1 \omega \text{ ALL } ( \text{ SELECT } \tau_2$   
 $\text{ FROM } \rho$   
 $\text{ WHERE } \varphi ) \Leftrightarrow$

NOT EXISTS ( SELECT  $\tau(\rho)$   
 FROM  $\rho$   
 WHERE  $(\varphi)$  AND  
 $\tau_1 \bar{\omega} \tau_2$  )

In law (I)  $\bar{\omega}$  is the negation of  $\omega$ , e.g.  $\bar{>} := >$ . In law (III) and (IV)  $\tau(\rho)$  refers to all attributes occurring in  $\rho$ . We assume  $\rho \equiv R_1 s_1, \dots, R_m s_m$ .

**Proof:** We assume  $\rho_1 \equiv s_1:R_1, \dots, s_{i-1}:R_{i-1}, s_{i+1}:R_{i+1}, \dots, s_m:R_m$ .

(I)  $sql2tc \llbracket \tau_1 \text{ IN ( SELECT } \tau_2$   
 FROM  $\rho$   
 WHERE  $\varphi$  )  $\rrbracket \Leftrightarrow$   
 $( \exists s_i:R_i ) ( ( \exists \rho_1 )$   
 $sql2tc \llbracket \varphi \rrbracket ) \wedge \tau_1 = \tau_2 ) \Leftrightarrow$   
 $sql2tc \llbracket \tau_1 = \text{ ANY ( SELECT } \tau_2$   
 FROM  $\rho$   
 WHERE  $\varphi$  )  $\rrbracket$

(II)  $sql2tc \llbracket \text{ NOT ( } \tau_1 \omega$   
 ALL ( SELECT  $\tau_2$   
 FROM  $\rho$   
 WHERE  $\varphi$  ) )  $\rrbracket \Leftrightarrow$   
 $\neg ( ( \forall s_i:R_i ) ( ( \exists \rho_1 )$   
 $sql2tc \llbracket \varphi \rrbracket ) \Rightarrow \tau_1 \omega \tau_2 ) \Leftrightarrow$   
 $\neg ( ( \forall s_i:R_i ) ( ( \forall \rho_1 )$   
 $\neg sql2tc \llbracket \varphi \rrbracket ) \vee \tau_1 \omega \tau_2 ) \Leftrightarrow$   
 $\neg ( ( \forall s_i:R_i ) ( \forall \rho_1 )$   
 $( \neg sql2tc \llbracket \varphi \rrbracket \vee \tau_1 \omega \tau_2 ) ) \Leftrightarrow$   
 $( \exists s_i:R_i ) ( \exists \rho_1 )$   
 $( sql2tc \llbracket \varphi \rrbracket \wedge \tau_1 \bar{\omega} \tau_2 ) \Leftrightarrow$   
 $sql2tc \llbracket \tau_1 \bar{\omega} \text{ ANY ( SELECT } \tau_2$   
 FROM  $\rho$   
 WHERE  $\varphi$  )  $\rrbracket$

(III) and (IV) can be proved analogously. q.e.d.

After looking at the equivalent, but longer translation from SQL to tuple calculus described in [PBG 89], we found it worthwhile to publish our translation. Indeed, it has already been presented in a course on databases at the Technical University Braunschweig in 1987 [Go 87].

The SQL subset considered here does not take into account aggregation (e.g. COUNT, AVG, etc.) nor grouping features (e.g. GROUP BY, HAVING, etc.). These SQL language features will be discussed in a forthcoming paper.

## References

- [Bü 87] G.von Bülzingsloewen: Translating and Optimizing SQL Queries Having Aggregates. Proc. 13th VLDB, Saratoga Press, 1987.
- [CG 85] S. Ceri, G. Gottlob: Translating SQL Into Relational Algebra: Optimization, Semantics, and Equivalence of SQL Queries. IEEE Trans. on SE, Vol. SE-11, No.4, 1985.
- [Da 87] C.J.Date: A Guide to the SQL Standard. Addison-Wesley, Reading, 1987.
- [Go 87] M.Gogolla: Begleitmaterial zur Vorlesung Datenbanksysteme. TU Braunschweig, Wintersemester 87/88, 1987.
- [Ma 83] D.Maier: The Theory of Databases. Computer Science Press, Rockville, 1983.
- [PBG 89] J.Paredaens, P.de Bra, M.Gyssens, D.van Gucht: The Structure of the Relational Database Model. EATCS Monographs on Theoretical Computer Science No. 17, Springer Verlag, Berlin, 1989.