

# Aquery to Q Compiler: Parser Grammar

Jose Cambronero

February 5, 2015

## 1 Introduction

As part of implementing a compiler from Aquery to q<sup>1</sup>, we have developed the BNF grammar below. This grammar is eventually used in implementing a flex/bison parser for Aquery.

## 2 Grammar

An Aquery program consists on a top level composed by queries, table/view creation, insert/update/delete statements. An empty program is also a valid aquery program as noted by the epsilon production for top-level.

### 2.1 Top-level program

$\langle program \rangle ::= \langle top\_level \rangle$

$\langle top\_level \rangle ::= \langle global\_query \rangle \langle top\_level \rangle$   
|  $\langle create\_table\_or\_view \rangle \langle top\_level \rangle$   
|  $\langle insert\_statement \rangle \langle top\_level \rangle$   
|  $\langle update\_statement \rangle \langle top\_level \rangle$   
|  $\langle delete\_statement \rangle \langle top\_level \rangle$   
|  $\langle user\_function\_definition \rangle \langle top\_level \rangle$   
|  $\epsilon$  (\*note: an empty aquery program is still an aquery program\*)

### 2.2 Local and global queries

We proceed to define what constitutes a global and local query, those that can solely be used within queries between the **WITH** keyword and the following global query. Note that aside from the necessary declarations at the beginning, the remainder of the query is a normal query, and thus refers to the grammar rule associated with the base\_query non-terminal.

$\langle global\_query \rangle ::= \langle local\_queries \rangle \langle base\_query \rangle$

---

<sup>1</sup>Aquery is an ordered database query language developed by Alberto Lerner and Dennis Shasha, for more information please see <https://cs.nyu.edu/web/Research/TechReports/TR2003-836/TR2003-836.pdf>

$$\begin{aligned} \langle local\_queries \rangle &::= 'WITH' \langle local\_query \rangle \langle local\_queries\_tail \rangle \\ &| \epsilon \\ \langle local\_queries\_tail \rangle &::= \langle local\_query \rangle \langle local\_queries\_tail \rangle \\ &| \epsilon \\ \langle local\_query \rangle &::= \langle identifier \rangle \langle col\_aliases \rangle 'AS' '(' \langle base\_query \rangle ')' \\ \langle col\_aliases \rangle &::= '(' \langle comma\_identifier\_list \rangle ')' \\ &| \epsilon \\ \langle comma\_identifier\_list \rangle &::= \langle identifier \rangle \langle comma\_identifier\_list\_tail \rangle \\ \langle comma\_identifier\_list\_tail \rangle &::= ',' \langle identifier \rangle \langle comma\_identifier\_list\_tail \rangle \\ &| \epsilon \\ \langle column\_list \rangle &::= \langle column\_elem \rangle \langle column\_list\_tail \rangle \\ \langle column\_elem \rangle &::= \langle identifier \rangle | \langle column\_access \rangle \\ \langle column\_access \rangle &::= \langle identifier \rangle '.' \langle identifier \rangle \\ \langle column\_list\_tail \rangle &::= ',' \langle column\_elem \rangle \langle column\_list\_tail \rangle \\ &| \epsilon \end{aligned}$$

### 2.3 Base query

A query requires a select clause and a from clause. There are additional optional clauses including an ordering clause (the base of declarative order in Aquery), a where clause, a group by clause, and a having clause, which can be used solely in conjunction with a group by clause.

$$\begin{aligned} \langle base\_query \rangle &::= \langle select\_clause \rangle \langle from\_clause \rangle \langle order\_clause \rangle \langle where\_clause \rangle \\ &\quad \langle groupby\_clause \rangle \\ \langle select\_clause \rangle &::= 'SELECT' \langle select\_elem \rangle \langle select\_clause\_tail \rangle \\ \langle select\_elem \rangle &::= \langle value\_expression \rangle 'as' \langle identifier \rangle \\ &| \langle value\_expression \rangle \\ \langle select\_clause\_tail \rangle &::= ',' \langle select\_elem \rangle \langle select\_clause\_tail \rangle \\ &| \epsilon \\ \langle from\_clause \rangle &::= 'FROM' \langle table\_expressions \rangle \\ \langle order\_clause \rangle &::= 'ASSUMING' \langle order\_specs \rangle \\ &| \epsilon \\ \langle order\_specs \rangle &::= \langle order\_spec \rangle \langle order\_specs\_tail \rangle \\ \langle order\_spec \rangle &::= 'ASC' \langle column\_list \rangle \\ &| 'DESC' \langle column\_list \rangle \end{aligned}$$

$\langle order\_specs\_tail \rangle ::= 'AND' \langle order\_spec \rangle \langle order\_specs\_tail \rangle$   
 $| \epsilon$

$\langle where\_clause \rangle ::= 'WHERE' \langle search\_condition \rangle$   
 $| \epsilon$

$\langle groupby\_clause \rangle ::= 'GROUP' 'BY' \langle comma\_value\_expression\_list \rangle \langle having\_clause \rangle$   
 $| \epsilon$

$\langle having\_clause \rangle ::= 'HAVING' \langle search\_condition \rangle$   
 $| \epsilon$

### 2.3.1 Search condition

**where** clauses, along with other clauses such as **having** and **on**, require a search condition non-terminal. We borrow the notion of a search condition from SQL 92 grammar, but also allow traditional value expressions to be a member of a search condition, which allows user defined functions to be used as stand-alone predicates, as well as boolean-typed columns. Please see <http://savage.net.au/SQL/sql-92.bnf.html#searchcondition> for the SQL92 details on this.

$\langle search\_condition \rangle ::= \langle boolean\_term \rangle$   
 $| \langle search\_condition \rangle 'OR' \langle boolean\_term \rangle$

$\langle boolean\_term \rangle ::= \langle boolean\_factor \rangle$   
 $| \langle boolean\_term \rangle 'AND' \langle boolean\_factor \rangle$

$\langle boolean\_factor \rangle ::= \langle boolean\_primary \rangle$   
 $| 'NOT' \langle boolean\_primary \rangle$

$\langle boolean\_primary \rangle ::= \langle predicate \rangle$   
 $| '(' \langle search\_condition \rangle ')'$

$\langle predicate \rangle ::= \langle value\_expression \rangle \langle postfix\_predicate \rangle$   
 $| \langle overlaps\_predicate \rangle$

$\langle postfix\_predicate \rangle ::= \langle between\_predicate \rangle$   
 $| \langle in\_predicate \rangle$   
 $| \langle like\_predicate \rangle$   
 $| \langle null\_predicate \rangle$   
 $| \langle is\_predicate \rangle$   
 $| \epsilon$  (\* value\_expression is already boolean \*)

$\langle between\_predicate \rangle ::= BETWEEN \langle value\_expression \rangle 'AND' \langle value\_expression \rangle$   
 $| 'NOT' 'BETWEEN' \langle value\_expression \rangle 'AND' \langle value\_expression \rangle$

$\langle in\_predicate \rangle ::= 'IN' \langle in\_pred\_spec \rangle$   
 $| 'NOT' 'IN' \langle in\_pred\_spec \rangle$

$\langle in\_pred\_spec \rangle ::= \langle value\_expression \rangle$  (\* implicit list \*)  
 $| '(' \langle comma\_value\_expression\_list \rangle ')'$

$\langle \text{like\_predicate} \rangle ::= \text{'LIKE'} \langle \text{value\_expression} \rangle$   
 $| \text{'NOT' 'LIKE'} \langle \text{value\_expression} \rangle$

$\langle \text{null\_predicate} \rangle ::= \text{'IS' 'NULL'}$   
 $| \text{'IS' 'NOT' 'NULL'}$

$\langle \text{is\_predicate} \rangle ::= \text{'IS'} \langle \text{truth\_value} \rangle$   
 $| \text{'IS' 'NOT'} \langle \text{truth\_value} \rangle$

$\langle \text{truth\_value} \rangle ::= \text{'TRUE'} | \text{'FALSE'}$

$\langle \text{overlaps\_predicate} \rangle ::= \langle \text{range\_value\_expression} \rangle \text{'OVERLAPS'} \langle \text{range\_value\_expression} \rangle$

$\langle \text{range\_value\_expression} \rangle ::= \text{'('} \langle \text{value\_expression} \rangle \text{' ,' } \langle \text{value\_expression} \rangle \text{'}'$

### 2.3.2 Table Expressions

We now proceed to define what table expressions constitute. Informally, table expression can be an identifier associated with a table, or an operation on a table (such as flatten), or a join on tables. Join grammar inspired by <http://savage.net.au/SQL/sql-92.bnf.html#joinedtable>. Note that the precedence/associativity of joins and other operations is directly encoded in the grammar.

$\langle \text{table\_expressions} \rangle ::= \langle \text{joined\_table} \rangle \langle \text{table\_expressions\_tail} \rangle$

$\langle \text{table\_expressions\_tail} \rangle ::= \text{' ,' } \langle \text{joined\_table} \rangle \langle \text{table\_expressions\_tail} \rangle$  (\*note  
the semantics of this are cross join\*)  
 $| \epsilon$

$\langle \text{joined\_table} \rangle ::= \langle \text{qualified\_join} \rangle$   
 $| \langle \text{qualified\_join} \rangle \text{'CROSS' 'JOIN'} \langle \text{joined\_table} \rangle$

$\langle \text{qualified\_join} \rangle ::= \langle \text{table\_expression} \rangle$   
 $| \langle \text{table\_expression} \rangle \text{'NATURAL' } \langle \text{join\_type} \rangle \text{'JOIN'} \langle \text{qualified\_join} \rangle$   
 $| \langle \text{table\_expression} \rangle \langle \text{join\_type} \rangle \text{'JOIN'} \langle \text{qualified\_join} \rangle \langle \text{join\_spec} \rangle$

$\langle \text{table\_expression} \rangle ::= \langle \text{table\_expression\_main} \rangle$   
 $| \langle \text{built\_in\_table\_fun} \rangle \text{'(' } \langle \text{table\_expression\_main} \rangle \text{'}'$

$\langle \text{table\_expression\_main} \rangle ::= \langle \text{identifier} \rangle \langle \text{identifier} \rangle$  (\* implicit correlation name  
\*)  
 $| \langle \text{identifier} \rangle \text{'AS' } \langle \text{identifier} \rangle$  (\* explicit correlation name \*)  
 $| \langle \text{identifier} \rangle$   
 $| \text{'(' } \langle \text{joined\_table} \rangle \text{'}'$

$\langle \text{built\_in\_table\_fun} \rangle ::= \text{'FLATTEN'}$  (\* potentially more to add here \*)

$\langle \text{join\_type} \rangle ::= \text{'INNER'} | \text{'LEFT' 'OUTER'} | \text{'LEFT' } | \text{'RIGHT' 'OUTER' } |$   
 $\text{'RIGHT' } | \text{'FULL' } | \text{'FULL' 'OUTER' } | \epsilon$

$\langle join\_spec \rangle ::= \langle on\_clause \rangle \mid \langle using\_clause \rangle$

$\langle on\_clause \rangle ::= 'ON' \langle search\_condition \rangle$

$\langle using\_clause \rangle ::= 'USING' '(' \langle comma\_identifier\_list \rangle ')'$

## 2.4 Table and View Creation

We define table and view creation at the top-level

$\langle create\_table\_or\_view \rangle ::= 'CREATE' 'TABLE' 'ID' \langle create\_spec \rangle$   
 $\mid 'CREATE' 'VIEW' 'ID' \langle create\_spec \rangle$

$\langle create\_spec \rangle ::= 'AS' \langle global\_query \rangle$   
 $\mid '(' \langle schema \rangle ')'$

$\langle schema \rangle ::= \langle schema\_element \rangle \langle schema\_tail \rangle$

$\langle schema\_element \rangle ::= \langle identifier \rangle \langle type \rangle$

$\langle schema\_tail \rangle ::= ',' \langle schema\_element \rangle \langle schema\_tail \rangle$   
 $\mid \epsilon$

$\langle type \rangle ::= 'INT' \mid 'FLOAT' \mid 'STRING' \mid 'DATE' \mid 'BOOLEAN' \mid 'HEX'$

## 2.5 Updating, Inserting, Deleting

We define the grammar relating to update, insert and delete statements at the top level.

$\langle update\_statement \rangle ::= 'UPDATE' \langle identifier \rangle 'SET' \langle set\_clauses \rangle \langle order\_clause \rangle$   
 $\langle where\_clause \rangle$

$\langle set\_clauses \rangle ::= \langle set\_clause \rangle \langle set\_clauses\_tail \rangle$

$\langle set\_clauses\_tail \rangle ::= ',' \langle set\_clause \rangle \langle set\_clauses\_tail \rangle$   
 $\mid \epsilon$

$\langle set\_clause \rangle ::= \langle identifier \rangle '=' \langle value\_expression \rangle$

$\langle insert\_statement \rangle ::= 'INSERT' 'INTO' \langle identifier \rangle \langle order\_clause \rangle \langle insert\_modifier \rangle$   
 $\langle insert\_source \rangle$

$\langle insert\_modifier \rangle ::= '(' \langle comma\_identifier\_list \rangle ')'$  (\* insert values into given  
order of column names \*)  
 $\mid \epsilon$  (\* insert into default column order \*)

$\langle insert\_source \rangle ::= \langle global\_query \rangle$   
 $\mid 'VALUES' '(' \langle comma\_value\_expression\_list \rangle ')'$

$\langle delete\_statement \rangle ::= 'DELETE' \langle from\_clause \rangle \langle order\_clause \rangle \langle where\_clause \rangle$   
 $| 'DELETE' \langle comma\_identifier\_list \rangle \langle from\_clause \rangle$  (\* similarly to q, we can  
 choose to delete rows where the predicates in where\_clause are true, or we  
 can choose to delete a column, but not both. No need to specify order,  
 since will delete whole column...\*)

## 2.6 User Defined Functions

We now define another element of the top-level: user defined function. Functions can have a series of expressions, queries, or local variable definitions. All but the last of which have to be followed by a semi-colon. The result of the function is the last expression evaluated.

$\langle user\_function\_definition \rangle ::= 'FUNCTION' \langle identifier \rangle '(' \langle comma\_identifier\_list \rangle$   
 $' ) ' \{ ' \langle function\_body \rangle ' \}$

$\langle function\_body \rangle ::= \langle function\_body\_elem \rangle \langle function\_body\_tail \rangle$   
 $| \epsilon$  (\*note: a function with no body is still a function \*)

$\langle function\_body\_tail \rangle ::= ';' \langle function\_body\_elem \rangle \langle function\_body\_tail \rangle$   
 $| \epsilon$

$\langle function\_body\_elem \rangle ::= \langle value\_expression \rangle | \langle function\_local\_var\_def \rangle | \langle local\_queries \rangle$   
 $\langle base\_query \rangle$  (\* functions can have expressions or queries \*)

$\langle function\_local\_var\_def \rangle ::= \langle identifier \rangle ':=' \langle value\_expression \rangle$

## 2.7 Value Expressions

We encode operator precedence and associativity into the grammar itself. This section of the grammar draws inspiration from <http://www.lysator.liu.se/c/ANSI-C-grammar-y.html>. Some expressions also draw inspiration from <http://savage.net.au/SQL/sql-92.bnf.html>.

$\langle constant \rangle ::= \langle integer \rangle | \langle float \rangle | \langle date \rangle | \langle string \rangle | \langle hex \rangle | \langle truth\_value \rangle$

$\langle table\_constant \rangle ::= 'ROWID' | \langle column\_access \rangle | '*'$

$\langle case\_expression \rangle ::= 'CASE' \langle case\_clause \rangle \langle when\_clauses \rangle \langle else\_clause \rangle 'END'$   
 (\* sql92 \*)

$\langle case\_clause \rangle ::= \langle value\_expression \rangle$  (\* value is compared to values in when  
 clauses\*)  
 $| \epsilon$  (\* when clause is treated as a search \*)

$\langle when\_clauses \rangle ::= \langle when\_clause \rangle \langle when\_clauses\_tail \rangle$

$\langle when\_clauses\_tail \rangle ::= \langle when\_clause \rangle \langle when\_clauses\_tail \rangle$   
 $| \epsilon$

$\langle \text{when\_clause} \rangle ::= \text{'WHEN'} \langle \text{search\_condition} \rangle \text{'THEN'} \langle \text{value\_expression} \rangle$  (\*search condition generalizes to the value-comparison case expression, since a search condition non-terminal can expand into a simple value expression \*)

$\langle \text{else\_clause} \rangle ::= \text{'ELSE'} \langle \text{value\_expression} \rangle$   
 $\quad | \quad \epsilon$

$\langle \text{main\_expression} \rangle ::= \langle \text{constant} \rangle | \langle \text{table\_constant} \rangle | \langle \text{identifier} \rangle | \text{'('} \langle \text{value\_expression} \rangle \text{'})' | \langle \text{case\_expression} \rangle$

$\langle \text{call} \rangle ::= \langle \text{main\_expression} \rangle$   
 $\quad | \quad \langle \text{main\_expression} \rangle \text{'['} \langle \text{indexing} \rangle \text{'}'}$   
 $\quad | \quad \langle \text{built\_in\_fun} \rangle \text{'('} \text{'}'}$   
 $\quad | \quad \langle \text{built\_in\_fun} \rangle \text{'('} \langle \text{comma\_value\_expression\_list} \rangle \text{'}'}$   
 $\quad | \quad \langle \text{identifier} \rangle \text{'('} \text{'}'}$  (\* user defined functions \*)  
 $\quad | \quad \langle \text{identifier} \rangle \text{'('} \langle \text{comma\_value\_expression\_list} \rangle \text{'}'}$  (\*user defined functions\*)

$\langle \text{indexing} \rangle ::= \text{ODD} | \text{EVEN} | \text{EVERY} \langle \text{integer} \rangle$

$\langle \text{built\_in\_fun} \rangle ::= \text{'abs'} | \text{'avg'} | \text{'count'} | \text{'deltas'} | \text{'distinct'} | \text{'drop'} | \text{'fill'} |$   
 $\quad \text{'first'} | \text{'last'} | \text{'max'} | \text{'maxs'} | \text{'min'} | \text{'mins'} | \text{'mod'} | \text{'next'} | \text{'not'} | \text{'prev'} |$   
 $\quad \text{'prd'} | \text{'prds'} | \text{'reverse'} | \text{'sum'} | \text{'sums'} | \text{'stddev'}$

$\langle \text{exp\_expression} \rangle ::= \langle \text{call} \rangle$   
 $\quad | \quad \langle \text{call} \rangle \langle \text{exp\_op} \rangle \langle \text{exp\_expression} \rangle$

$\langle \text{mult\_expression} \rangle ::= \langle \text{exp\_expression} \rangle$   
 $\quad | \quad \langle \text{mult\_expression} \rangle \langle \text{times\_op} \rangle \langle \text{exp\_expression} \rangle$   
 $\quad | \quad \langle \text{mult\_expression} \rangle \langle \text{div\_op} \rangle \langle \text{exp\_expression} \rangle$

$\langle \text{add\_expression} \rangle ::= \langle \text{mult\_expression} \rangle$   
 $\quad | \quad \langle \text{add\_expression} \rangle \langle \text{plus\_op} \rangle \langle \text{mult\_expression} \rangle$   
 $\quad | \quad \langle \text{add\_expression} \rangle \langle \text{minus\_op} \rangle \langle \text{mult\_expression} \rangle$

$\langle \text{rel\_expression} \rangle ::= \langle \text{add\_expression} \rangle$   
 $\quad | \quad \langle \text{rel\_expression} \rangle \langle \text{l\_op} \rangle \langle \text{add\_expression} \rangle$   
 $\quad | \quad \langle \text{rel\_expression} \rangle \langle \text{g\_op} \rangle \langle \text{add\_expression} \rangle$   
 $\quad | \quad \langle \text{rel\_expression} \rangle \langle \text{le\_op} \rangle \langle \text{add\_expression} \rangle$   
 $\quad | \quad \langle \text{rel\_expression} \rangle \langle \text{ge\_op} \rangle \langle \text{add\_expression} \rangle$

$\langle \text{eq\_expression} \rangle ::= \langle \text{rel\_expression} \rangle$   
 $\quad | \quad \langle \text{eq\_expression} \rangle \langle \text{eq\_op} \rangle \langle \text{rel\_expression} \rangle$   
 $\quad | \quad \langle \text{eq\_expression} \rangle \langle \text{neq\_op} \rangle \langle \text{rel\_expression} \rangle$

$\langle \text{and\_expression} \rangle ::= \langle \text{eq\_expression} \rangle$   
 $\quad | \quad \langle \text{and\_expression} \rangle \langle \text{and\_op} \rangle \langle \text{eq\_expression} \rangle$

$\langle \text{or\_expression} \rangle ::= \langle \text{and\_expression} \rangle$   
 $\quad | \quad \langle \text{or\_expression} \rangle \langle \text{or\_op} \rangle \langle \text{and\_expression} \rangle$

$\langle \text{value\_expression} \rangle ::= \langle \text{or\_expression} \rangle$

Now that we have value expressions defined, we define a form of value expression list: comma separated value expressions.

$$\langle \text{comma\_value\_expression\_list} \rangle ::= \langle \text{value\_expression} \rangle \langle \text{comma\_value\_expression\_list\_tail} \rangle$$
$$\langle \text{comma\_value\_expression\_list\_tail} \rangle ::= ', ' \langle \text{value\_expression} \rangle \langle \text{comma\_value\_expression\_list\_tail} \rangle$$
$$| \epsilon$$

This concludes the formal outline of the Aquery grammar. Note that this grammar maybe revised and changed as necessary throughout development if need be.

For a flex/bison implementation of this grammar please see <https://www.github.com/josepablocam/aquery2q/parser/>