

Aquery to Q Compiler: Parser Grammar

Jose Cambronero

January 21, 2015

1 Introduction

As part of implementing a compiler from Aquery to q¹, 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 of a list of semi-colon separated queries with global context. Global queries in turn consist of a potential list of local queries followed by a global query.

$$\langle program \rangle ::= \langle global_queries \rangle$$
$$\begin{aligned} \langle global_queries \rangle &::= \langle global_query \rangle \langle global_queries \rangle \\ &\quad | \quad \epsilon \end{aligned}$$
$$\langle global_query \rangle ::= \langle local_queries \rangle \langle query \rangle \text{';'}$$

We proceed to define what constitutes a 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 query non-terminal.

$$\langle comma_identifier_list \rangle ::= \langle identifier \rangle \langle comma_identifier_list_tail \rangle$$
$$\begin{aligned} \langle comma_identifier_list_tail \rangle &::= \text{' ,' } \langle identifier \rangle \langle comma_identifier_list_tail \rangle \\ &\quad | \quad \epsilon \end{aligned}$$
$$\begin{aligned} \langle local_queries \rangle &::= \text{'WITH' } \langle local_query \rangle \langle local_queries_tail \rangle \\ &\quad | \quad \epsilon \end{aligned}$$
$$\begin{aligned} \langle local_queries_tail \rangle &::= \langle local_query \rangle \langle local_queries_tail \rangle \\ &\quad | \quad \epsilon \end{aligned}$$

¹Aquery 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>

$$\langle local_query \rangle ::= \langle identifier \rangle \langle col_aliases \rangle 'AS' '(' \langle query \rangle ')'$$

$$\langle col_aliases \rangle ::= '(' \langle comma_identifier_list \rangle ')'$$

$$| \epsilon$$

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, and a group by clause

$$\langle query \rangle ::= \langle select_clause \rangle \langle from_clause \rangle \langle order_clause \rangle \langle where_clause \rangle \langle groupby_clause \rangle$$

We breakout the grammar for each relevant clause below

$$\langle select_clause \rangle ::= \langle select_elem \rangle \langle select_clause_tail \rangle$$

$$\langle select_elem \rangle ::= \langle expression \rangle 'as' \langle identifier \rangle$$

$$| \langle expression \rangle$$

$$\langle select_clause_tail \rangle ::= ', ' \langle select_elem \rangle \langle select_clause_tail \rangle$$

$$| \epsilon$$

$$\langle from_clause \rangle ::= 'FROM' \langle table_expressions \rangle$$

$$\langle order_clause \rangle ::= 'ASSUMING' 'ORDER' \langle comma_identifier_list \rangle (*comment:$$

does the order clause allow expression on the fly, similarly to group by? *)

$$| \epsilon$$

$$\langle where_clause \rangle ::= 'WHERE' \langle and_expression_list \rangle$$

$$| \epsilon$$

$$\langle groupby_clause \rangle ::= 'GROUP' 'BY' \langle comma_expression_list \rangle$$

$$| \epsilon$$

We now proceed to define what table expressions constitute. Table expression can be an identifier associated with a table, or an operation on a table (such as flatten).

$$\langle table_expressions \rangle ::= \langle table_expression \rangle \langle table_expression_tail \rangle$$

$$\langle table_expression_tail \rangle ::= ', ' \langle table_expression \rangle \langle table_expression_tail \rangle$$

$$| \epsilon$$

$$\langle table_expression \rangle ::= \langle table_exp \rangle \langle identifier \rangle$$

$$| \langle table_exp \rangle$$

$$\langle table_exp \rangle ::= 'FLATTEN' '(' \langle identifier \rangle ')'$$

$$| \langle identifier \rangle$$

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>.

$\langle literal \rangle ::= \langle identifier \rangle \mid '*' \mid \langle column_access \rangle \mid \langle integer \rangle \mid \langle float \rangle \mid \langle date \rangle \mid \langle string \rangle \mid '(' \langle expression \rangle ')'$

$\langle column_access \rangle ::= \langle identifier \rangle '.' \langle identifier \rangle$

$\langle call \rangle ::= \langle literal \rangle$
 $\mid \langle literal \rangle '[' \langle indexing \rangle ']'$
 $\mid \langle built_in \rangle '(' ' ')$
 $\mid \langle built_in \rangle '(' \langle comma_expression_list \rangle ')'$

$\langle indexing \rangle ::= ODD \mid EVEN \mid EVERY \langle integer \rangle$

$\langle built_in \rangle ::= 'abs' \mid 'avg' \mid 'count' \mid 'deltas' \mid 'distinct' \mid 'drop' \mid 'first' \mid 'last'$
 $\mid 'max' \mid 'maxs' \mid 'min' \mid 'mins' \mid 'mod' \mid 'next' \mid 'prev' \mid 'prd' \mid 'prds'$
 $\mid 'reverse' \mid 'sum' \mid 'sums' \mid 'stddev'$

$\langle mult_expression \rangle ::= \langle call \rangle$
 $\mid \langle mult_expression \rangle \langle times_op \rangle \langle call \rangle$
 $\mid \langle mult_expression \rangle \langle div_op \rangle \langle call \rangle$

$\langle add_expression \rangle ::= \langle mult_expression \rangle$
 $\mid \langle add_expression \rangle \langle plus_op \rangle \langle mult_expression \rangle$
 $\mid \langle add_expression \rangle \langle minus_op \rangle \langle mult_expression \rangle$

$\langle rel_expression \rangle ::= \langle add_expression \rangle$
 $\mid \langle rel_expression \rangle \langle l_op \rangle \langle add_expression \rangle$
 $\mid \langle rel_expression \rangle \langle g_op \rangle \langle add_expression \rangle$
 $\mid \langle rel_expression \rangle \langle le_op \rangle \langle add_expression \rangle$
 $\mid \langle rel_expression \rangle \langle ge_op \rangle \langle add_expression \rangle$

$\langle eq_expression \rangle ::= \langle rel_expression \rangle$
 $\mid \langle eq_expression \rangle \langle eq_op \rangle \langle rel_expression \rangle$
 $\mid \langle eq_expression \rangle \langle neq_op \rangle \langle rel_expression \rangle$

$\langle and_expression \rangle ::= \langle eq_expression \rangle$
 $\mid \langle and_expression \rangle \langle and_op \rangle \langle eq_expression \rangle$

$\langle or_expression \rangle ::= \langle and_expression \rangle$
 $\mid \langle or_expression \rangle ::= \langle or_expression \rangle \langle or_op \rangle \langle and_expression \rangle$

$\langle logical_expression \rangle ::= \langle and_expression \rangle \mid \langle or_expression \rangle$

$\langle expression \rangle ::= \langle logical_expression \rangle$

Now that we have expressions defined, we define 2 forms of expression lists, comma and **AND** separated expression lists.

$\langle comma_expression_list \rangle ::= \langle expression \rangle \langle comma_expression_list_tail \rangle$

$\langle comma_expression_list_tail \rangle ::= ', ' \langle expression \rangle \langle comma_expression_list_tail \rangle$
 $\mid \epsilon$

$$\langle and_expression_list \rangle ::= \langle expression \rangle \langle and_expression_list_tail \rangle$$
$$\begin{aligned} \langle and_expression_list_tail \rangle &::= 'AND' \langle expression \rangle \langle and_expression_list_tail \rangle \\ &| \epsilon \end{aligned}$$

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/>