

## 1. Introduction

These are the formal definitions of an Extended Relational Algebra in a style similar to Algebra A in Appendix A of DTATRM by D&D. The operators of this ERA are used to construct shorthands, intended to form the basis of a relational programming language. It should be noted that Appendix A takes the reduction path from language shorthand to algebra operation. Here the path is taken in the opposite direction, defining algebraic operations from which the shorthands of a language may be constructed.

These operators do not provide or require any particular type system. Any use of the algebra is dependent on an external provider of a type system and functions over those types. Shorthands such as `GROUP` and `UNGROUP`, `WRAP` and `UNWRAP` require that relations and tuples be provided by that type system.

The naming conventions are adapted from Appendix A. In what follows:

- A is an attribute name
- T is a type
- v is a value
- $\langle A, T \rangle$  is an attribute (member of a heading)
- $\langle A, T, v \rangle$  is an attribute value (member of a tuple)
- r is a relation
- Hr is the heading of r
- Br is the body of r
- tr is a tuple of r

## 2. Basics from Appendix A

These are the basic 5 definitions adapted from Appendix A.

### a) NOT

Given a relation  $r$  then  $s = \mathbf{NOT}(r)$  is defined as follows.

$$Hs = Hr$$
$$Bs = \{ ts : \exists tr \notin Br, ts = tr \}$$

The **NOT** operator yields the complement  $s$  of a given relation  $r$ . The heading of  $s$  is the heading of  $r$ . The body of  $s$  contains every tuple with that heading that is not in the body of  $r$ .

The **NOT** operator is used in shorthands that involve negation, such as `ANTIJOIN`, `MINUS`.

### b) REMOVE

Given a relation  $r$  and a type  $T$  such that  $\langle A, T \rangle \in Hr$  then  $s = \mathbf{REMOVE}(A)$  is defined as follows.

$$Hs = Hr \text{ minus } \{ \langle A, T \rangle \}$$
$$Bs = \{ ts : \exists tr \notin Br, \exists v \in T, \langle A, T, v \rangle \in tr, ts = tr \text{ minus } \{ \langle A, T, v \rangle \} \}$$

The **REMOVE** operator yields a relation  $s$  formed by removing a given attribute  $A$  from a given relation  $r$ . The operation is equivalent to taking the projection of  $r$  over all of its attributes except  $A$ . The heading of  $s$  is the heading of  $r$  minus the ordered pair  $\langle A, T \rangle$ . The body of  $s$  contains every tuple that conforms to the heading of  $s$  and is a subset of some tuple of  $r$ .

The **REMOVE** operator is used in shorthands that remove one or more attributes, such as `PROJECT`.

### c) RENAME

Given a relation  $r$ , a type  $T$  such that  $\langle A, T \rangle \in Hr$  and no type  $T$  such that  $\langle B, T \rangle \in Hr$  then  $s = \text{RENAME}(A, B)$  is defined as follows.

$$\begin{aligned} Hs &= (Hr \text{ minus } \{ \langle A, T \rangle \} ) \text{ union } \{ \langle B, T \rangle \} \\ Bs &= \{ ts : \exists tr \in Br, \exists v \in T, \langle A, T, v \rangle \in tr, \\ &\quad ts = (tr \text{ minus } \{ \langle A, T, v \rangle \} ) \text{ union } \{ \langle B, T, v \rangle \} \} \end{aligned}$$

The **RENAME** operator yields a relation  $s$  that differs from a given relation  $r$  only in the name of one of its attributes, which is changed from  $A$  to  $B$ . The heading of  $s$  is the heading of  $r$  except that the ordered pair  $\langle A, T \rangle$  is replaced by the ordered pair  $\langle B, T \rangle$ . The body of  $s$  consists of every tuple of the body of  $r$ , except that in each such tuple the triple  $\langle A, T, v \rangle$  is replaced by the triple  $\langle B, T, v \rangle$ .

The **RENAME** operator is used in shorthands that rename one or more attributes, such as **RENAME**.

### d) AND

Given relations  $r_1$  and  $r_2$  such that  $\langle A, T_1 \rangle \in Hr_1$  and  $\langle A, T_2 \rangle \in Hr_2$  implies  $T_1 = T_2$  then  $s = \text{AND}(r_1, r_2)$  is defined as follows.

$$\begin{aligned} Hs &= Hr_1 \text{ union } Hr_2 \\ Bs &= \{ ts : \exists tr_1 \in Br_1, \exists tr_2 \in Br_2, ts = tr_1 \text{ union } tr_2 \} \end{aligned}$$

The **AND** operator is relational conjunction, which is usually known as a natural join. The heading of  $s$  is the union of the headings of  $r_1$  and  $r_2$ . The body of  $s$  contains every tuple that conforms to the heading of  $s$  and is a superset of both some tuple in the body of  $r_1$  and some tuple in the body of  $r_2$ .

The **AND** operator is used in join shorthands such as **JOIN**, **SEMIJOIN**, **ANTIJOIN**, **COMPOSE**, **TIMES**.

### e) OR

Given relations  $r_1$  and  $r_2$  such that  $\langle A, T_1 \rangle \in Hr_1$  and  $\langle A, T_2 \rangle \in Hr_2$  implies  $T_1 = T_2$  then  $s = \text{OR}(r_1, r_2)$  is defined as follows.

$$\begin{aligned} Hs &= Hr_1 \text{ union } Hr_2 \\ Bs &= \{ ts : \exists tr_1, \exists tr_2, (tr_1 \in Br_1 \text{ or } tr_2 \in Br_2), ts = tr_1 \text{ union } tr_2 \} \end{aligned}$$

The **OR** operator is relational disjunction, being a generalization of what is usually known as union. In the special case where the given relations  $r_1$  and  $r_2$  have the same heading, the result  $s$  is in fact the union of those two relations in the traditional sense. The heading of  $s$  is the union of the headings of  $r_1$  and  $r_2$ . The body of  $s$  contains every tuple that conforms to the heading of  $s$  and is a superset of either some tuple in the body of  $r_1$  or some tuple in the body of  $r_2$ .

The **OR** operator is used in set shorthands such as **UNION**, **MINUS**, **SYMDIFF**, **INSERT**.

## 3. Relational Constant

This is the formalisation of a relational constant or 'relcon', and relies on some previously defined function  $f$ . Separate formalisations are provided for monadic and dyadic functions, in both predicate and value-returning forms. Extending them to  $n$ -adic functions is left as an exercise.

### a) Monadic predicate

Given attribute names  $X$ , types  $T_x$  and function  $f$  with the type signature  $T_x \rightarrow \text{Boolean}$  then  $s = \text{RELCON}(X, f)$  is defined as follows.

$$Hs = \{ \langle X, T_x \rangle \}$$

$$Bs = \{ ts : \exists vx \in Tx, f(vx) = true, ts = \{ \langle X, Tx, vx \rangle \} \}$$

#### b) Monadic value

Given attribute names  $X, Y$ , types  $T_x, T_y$  and function  $f$  with the type signature  $T_x \rightarrow T_y$  then  $s = \mathbf{RELCON}(X, Y, f)$  is defined as follows.

$$Hs = \{ \langle X, T_x \rangle, \langle Y, T_y \rangle \}$$

$$Bs = \{ ts : \exists vx \in T_x, \exists vy \in T_y, f(vx) = vy, \\ ts = \{ \langle X, T_x, vx \rangle, \langle Y, T_y, vy \rangle \} \}$$

#### c) Dyadic predicate

Given attribute names  $X, Y$ , types  $T_x, T_y$  and function  $f$  with the type signature  $T_x \rightarrow T_y \rightarrow \text{Boolean}$  then  $s = \mathbf{RELCON}(X, Y, f)$  is defined as follows.

$$Hs = \{ \langle X, T_x \rangle, \langle Y, T_y \rangle \}$$

$$Bs = \{ ts : \exists vx \in T_x, \exists vy \in T_y, f(vx, vy) = true, \\ ts = \{ \langle X, T_x, vx \rangle, \langle Y, T_y, vy \rangle \} \}$$

#### d) Dyadic value

Given attribute names  $X, Y, Z$ , types  $T_x, T_y, T_z$  and function  $f$  with the type signature  $T_x \rightarrow T_y \rightarrow T_z$  then  $s = \mathbf{RELCON}(X, Y, Z, f)$  is defined as follows.

$$Hs = \{ \langle X, T_x \rangle, \langle Y, T_y \rangle, \langle Z, T_z \rangle \}$$

$$Bs = \{ ts : \exists vx \in T_x, \exists vy \in T_y, \exists vz \in T_z, f(vx, vy) = vz, \\ ts = \{ \langle X, T_x, vx \rangle, \langle Y, T_y, vy \rangle, \langle Z, T_z, vz \rangle \} \}$$

[Note: for the **relcon PLUS** in App-A, types  $T_x, T_y$  and  $T_z$  are all **INTEGER**, and the function  $f$  is the scalar operator "+".]

The **RELCON** operator predicate forms are used together with **AND** in shorthands such as **WHERE/RESTRICT**.

The **RELCON** operator value-returning forms are used together with **AND** in shorthands such as **EXTEND, SUMMARIZE, UPDATE** and **DELETE**.

## 4. Aggregation

This is the formalisation of simple aggregation, similar to **SUM, MAX** or **MIN** but allowing for a wider range of functions. It relies on some previously defined aggregating function  $f$ . An aggregating function has an argument which is a list of values (a 'bag' or 'multiset' rather than a 'set') and returns a single value. Extending this approach to more complex aggregations is left as an exercise.

Given a relation  $r$ , a type  $T$  such that  $\langle A, T \rangle \in H_r$  and an aggregating function  $f$  with the type signature  $T[] \rightarrow T_a$  then  $s = \mathbf{AGGREGATE}(r, A, f)$  is defined as follows.

$$Hs = ( H_r \text{ minus } \{ \langle A, T \rangle \} ) \text{ union } \{ \langle A, T_a \rangle \}$$

$$Bs = \{ ts : \exists v \in T, \exists tr \in r, \exists \langle A, T, v \rangle \in tr, v \in v[], \\ ts = tr \text{ minus } \{ \langle A, T, v \rangle \} \text{ union } \{ \langle A, T_a, f(v[]) \rangle \} \}$$

Note: the term  $v[]$  should be read as: a list of the values of  $v$  for some value of  $ts$ .

The **AGGREGATION** operator is used in shorthands such as **SUMMARIZE**.

## 5. Transitive Closure

As transitive closure cannot be defined directly in set-builder notation, this formalisation defines it indirectly by means of a recurrence relation. This consists of a starting value (given) and a sequence of successors (defined by set-builder). The transitive closure is the fix-point union of that sequence.

The starting value ('seed') represents known edges in a directed graph; the end value is all the possible paths through the graph.

### a) Simple Transitive Closure

Given a relation  $r$  with the heading  $\{ \langle A, T \rangle, \langle B, T \rangle \}$  for some type  $T$  then the successor relation  $r'$  is defined as follows.

$$\begin{aligned} Hr' &= Hr \\ Br' &= \{ tr' : tr' \in Br \text{ or} \\ &\quad \exists v1 \in T, \exists v2 \in T, \exists v3 \in T, \exists v4 \in T, \\ &\quad \exists tr1 \in Br, tr1 = \{ \langle A, T, v1 \rangle, \langle B, T, v2 \rangle \}, \\ &\quad \exists tr2 \in Br, tr2 = \{ \langle A, T, v2 \rangle, \langle B, T, v3 \rangle \}, \\ &\quad tr' = \{ \langle A, T, v1 \rangle, \langle B, T, v3 \rangle \} \} \end{aligned}$$

Then the transitive closure  $s = \mathbf{TCLOSE}(r)$  is defined as follows.

$$S = r^1 \cup r^2 \cup r^3 \dots r^\infty$$

This is a linear recurrence, which can be shown to reach a fix-point termination.

The **TCLOSE** operator is used in a shorthand of the same name.

### b) Extended Transitive Closure

In this case each tuple is associated with a value, and this definition relies on some previously defined dyadic function  $f$  that takes values of that type as its argument and return value. The value computed represents in some sense the cost of that path.

Given a relation  $r$  with heading  $\{ \langle A, T \rangle, \langle B, T \rangle, \langle C, Tv \rangle \}$  for some types  $T$  and  $Tv$ , and a dyadic function  $f$  with the type signature  $Tv \rightarrow Tv \rightarrow Tv$  then the successor relation  $r'$  is defined as follows.

$$\begin{aligned} Hr' &= Hr \\ Br' &= \{ tr' : tr' \in Br \text{ or} \\ &\quad ( \exists v1 \in T, \exists v2 \in T, \exists v3 \in T, \exists v4 \in T, \\ &\quad \exists w1 \in Tv, \exists w2 \in Tv, \\ &\quad \exists tr1 \in Br, tr1 = \{ \langle A, T, v1 \rangle, \langle B, T, v2 \rangle, \langle C, Tv, w1 \rangle \}, \\ &\quad \exists tr2 \in Br, tr2 = \{ \langle A, T, v2 \rangle, \langle B, T, v3 \rangle, \langle C, Tv, w2 \rangle \}, \\ &\quad tr' = \{ \langle A, T, v1 \rangle, \langle B, T, v3 \rangle, \langle C, Tv, f(w1, w2) \rangle \} ) \} \end{aligned}$$

Then the extended transitive closure  $s = \mathbf{ETCLOSE}(r, f)$  is defined as follows.

$$s = r^1 \cup r^2 \cup r^3 \dots r^\infty$$

This is a linear recurrence, which can be shown to reach a fix-point termination.

The **ETCLOSE** operator is used in combination with **AGGREGATE** in a shorthand to perform Generalised Transitive Closure as defined by D&D (RM VSS 5).