

Formalisation of Relational Algebra and a SQL-like Language with the RISCAL Model Checker

Joachim Borya

Johannes Kepler Universität

January 26th, 2023

Outline

- 1 Topics
- 2 Recap on Relational Algebra
- 3 Implementation in RISCAL
- 4 Model Checking Algorithms and Theorems
- 5 A SQL-like wrapper language
- 6 Translation into standardized SQL
- 7 Current progress and further goals

- Formalization and verification of Relational Algebra with RISCAL
 - implementation of the common operations
 - model checking for proving correctness and theorems regarding the Relational Algebra
- Creation of a SQL-like language within RISCAL
- Demonstration by various examples

Recap on Relational Algebra

Definition

Let $n \in \mathbb{N}$ and $\mathcal{A} = \{A_1, \dots, A_n\}$ be an attribute schema. Then the set

$$\text{Row}_{\mathcal{A}} := \left\{ t : A \rightarrow \bigcup_{i=1}^n \text{dom}(A_i) : t(A_i) \in \text{dom}(A_i) \text{ for all } 1 \leq i \leq n \right\}$$

contains all tuples. Furthermore the set

$$\text{Relation}_{\mathcal{A}} := \{r \subseteq \text{Row}_{\mathcal{A}} : |r| < \infty\}$$

contains all relations.

- Attributes are identifiers for columns. They are combined to get an attribute schema.
- Any attribute has a set of allowed values, a so called domain.
- The algebra consists of these types and some operations.

Recap on Relational Algebra

```
val M:N; // maximum cardinality of relations
val N:N; // maximum length of rows/tuples
val K:N; // maximum number of tables
val D:N; // maximum query depth
val L:N; // maximum list length

type Element = N[1];
type Attribute = N[N-1];
type Length = N[N];
type TableId = N[K];
type Row = Map[Attribute, Element];
type Relation = Record[len:Length, tup:Set[Row]]
with |value.tup| ≤ M ∧ ∀ t:Row, i:Attribute. t ∈ value.tup ∧ i ≥ value.len ⇒ t[i] = 0;
type Database = Array[K, Relation];
```

Recap on Relational Algebra

- **Set operations**
- **Cartesian product**
- **Join:** Linking tables based on attributes
- **Selection:** Filtering on rows
- **Projection:** Filtering on columns

Implementation in RISCAL

Given an operator to implement:

- ① Start off with a predicate representing the specification of an operation.
- ② Rapid Prototyping: Choose a result that conforms to the predicate.
- ③ Algorithm: Implement it in terms of ...
 - RISCAL's high-level features like set constructors
 - and (only) loops, conditionals and variable assignments.

After the operations are verified we can check several theorems involving them.

Implementation in RISCAL

Definition

Let $\mathcal{A} = \{A_1, \dots, A_n\}$, $\mathcal{B} = \{B_1, \dots, B_m\}$ be relation schemas. The function $c : \text{Row}_{\mathcal{A}} \times \text{Row}_{\mathcal{B}} \rightarrow \text{Row}_{\mathcal{A} \cup \mathcal{B}}$ given by

$$c(t, u) := (t(A_1), \dots, t(A_n), u(B_1), \dots, u(B_m))$$

is called the concatenation of the rows t and u .

RISCAL

```
pred concat_spec(t:Row, t1:Row, t2:Row, n1:Length, n2:Length) ⇔
∀ i:Attribute. (
  if i < n1 then t[i] = t1[i]
  else if i ≥ n1 ∧ i < n1+n2 then t[i] = t2[i-n1]
  else t[i] = 0
);
```

Implementation in RISCAL

```
fun concat1(t1:Row, t2:Row, n1:Length, n2:Length):Row
requires n1 + n2 ≤ N;
= choose t:Row with concat_spec(t,t1,t2,n1,n2);
```

```
proc concat2(t1:Row, t2:Row, n1:Length, n2:Length):Row
requires n1 + n2 ≤ N;
ensures concat_spec(result,t1,t2,n1,n2); {
  var t:Row := Array[N,Element](0);
  for var i:Length:=0; i<n1; i:=i+1 do {
    t[i] := t1[i];
  }
  for var i:Length:=n1; i<n1+n2; i:=i+1 do {
    t[i] := t2[i-n1];
  }
  return t;
}
```

Definition

Let $\mathcal{A} = \{A_1, \dots, A_n\}$, $\mathcal{B} = \{B_1, \dots, B_m\}$ be attribute schemas. The cartesian product $C : \text{Relation}_{\mathcal{A}} \times \text{Relation}_{\mathcal{B}} \rightarrow \text{Relation}_{\mathcal{A} \cup \mathcal{B}}$ is given by

$$C(r, s) := \{c(t, u) : t \in r \wedge u \in s\}.$$

RISCAL

```
pred cartesian_spec(r:Relation, r1:Relation, r2:Relation) ⇔  
  r.len = r1.len+r2.len ∧  
  ∀ t:Row. t∈r.tup ⇔ ∃ t1:Row, t2:Row.  
    t1∈r1.tup ∧ t2∈r2.tup ∧ t = concat1(t1,t2,r1.len,r2.len);
```

Implementation in RISCAL

```
fun cartesian1(r1:Relation, r2:Relation):Relation
requires r1.len+r2.len ≤ N ∧ |r1.tup|*|r2.tup| ≤ M;
  = choose r:Relation with cartesian_spec(r,r1,r2);
```

```
fun cartesian2(r1:Relation, r2:Relation):Relation
requires r1.len+r2.len ≤ N ∧ |r1.tup|*|r2.tup| ≤ M;
ensures cartesian_spec(result,r1,r2);
  = ⟨len: r1.len+r2.len, tup: {concat1(t1,t2,r1.len,r2.len) | t1∈r1.tup, t2∈r2.tup}⟩;
```

```
proc cartesian3(r1:Relation, r2:Relation):Relation
requires r1.len+r2.len ≤ N ∧ |r1.tup|*|r2.tup| ≤ M;
ensures cartesian_spec(result,r1,r2); {
  var q:Relation := ⟨len: r1.len+r2.len, tup: choose s:Set[Row] with |s|=0⟩;

  for t1 ∈ r1.tup do {
    for t2 ∈ r2.tup do {
      q.tup := q.tup ∪ {concat2(t1, t2, r1.len, r2.len)};
    }
  }

  return q;
}
```

Implementation in RISCAL

Definition

Let $\mathcal{A} = \{A_1, \dots, A_n\}$, $\mathcal{B} = \{B_1, \dots, B_m\}$ be attribute schemas. The function $\bowtie: \text{Relation}_{\mathcal{A}} \times \text{Relation}_{\mathcal{B}} \times \mathcal{A} \times \mathcal{B} \rightarrow \text{Relation}_{\mathcal{A} \cup \mathcal{B}}$ given by

$$\bowtie(r, s, A, B) := \{c(t, u) : t \in r \wedge u \in s \wedge t(A) = u(B)\}$$

is called the (equi-)join of r and s on A and B .

RISCAL

```
pred join_spec(s:Relation, r1:Relation, r2:Relation, n1:Attribute, n2:Attribute) ⇔
s.len = r1.len+r2.len ∧
∀ t:Row. t∈s.tup ⇔ ∃ t1:Row, t2:Row.
(t1∈r1.tup ∧ t2∈r2.tup ∧ t = concat1(t1,t2,r1.len,r2.len) ∧ t1[n1] = t2[n2]);
```

Implementation in RISCAL

```
fun join1(r1:Relation, r2:Relation, n1:Attribute, n2:Attribute):Relation
requires n1<r1.len ^ n2<r2.len ^ r1.len+r2.len ≤ N ^ |r1.tup|*|r2.tup| ≤ M;
= choose s:Relation with join_spec(s,r1,r2,n1,n2);
```

```
fun join2(r1:Relation, r2:Relation, n1:Attribute, n2:Attribute):Relation
requires n1<r1.len ^ n2<r2.len ^ r1.len+r2.len ≤ N ^ |r1.tup|*|r2.tup| ≤ M;
ensures join_spec(result,r1,r2,n1,n2);
= {len: r1.len+r2.len, tup:
  {concat1(t1,t2,r1.len,r2.len) | t1∈r1.tup, t2∈r2.tup with t1[n1] = t2[n2]}};
```

```
proc join3(r1:Relation, r2:Relation, n1:Attribute, n2:Attribute):Relation
requires n1<r1.len ^ n2<r2.len ^ r1.len+r2.len ≤ N ^ |r1.tup|*|r2.tup| ≤ M;
ensures join_spec(result,r1,r2,n1,n2); {
  var q:Relation := {len: r1.len+r2.len, tup: choose s:Set[Row] with |s|=0};

  for t1 ∈ r1.tup do {
    for t2 ∈ r2.tup do {
      if t1[n1] = t2[n2] then {
        q.tup := q.tup ∪ {concat2(t1, t2, r1.len, r2.len)};
      }
    }
  }

  return q;
}
```

Implementation in RISCAL

Definition

Let $\mathcal{A} = \{A_1, \dots, A_n\}$ be an attribute schema. Then the function $\sigma : \text{Relation}_{\mathcal{A}} \times \mathcal{A} \times \bigcup_{A \in \mathcal{A}} \text{dom}(A) \rightarrow \text{Relation}_{\mathcal{A}}$ given by

$$\sigma(r, A, a) := \{t \in r : t(A) = a\}$$

is called the selection of all rows t in r with value a at attribute A .

RISCAL

```
pred select_spec(s:Relation, r:Relation, a:Attribute, e:Element) ⇔  
s.len = r.len ∧ ∀ t:Row. t∈s.tup ⇔ t∈r.tup ∧ t[a] = e;
```

Implementation in RISCAL

```
fun select1(r:Relation, a:Attribute, e:Element):Relation
requires a < r.len;
= choose s:Relation with select_spec(s,r,a,e);
```

```
fun select2(r:Relation, a:Attribute, e:Element):Relation
requires a < r.len;
ensures select_spec(result,r,a,e);
= ⟨len: r.len, tup: {t | t∈r.tup with t[a] = e}⟩;
```

```
proc select3(r:Relation, a:Attribute, e:Element):Relation
requires a < r.len;
ensures select_spec(result,r,a,e); {
  var q:Relation := ⟨len: r.len, tup: choose s:Set[Row] with |s|=0⟩;

  for t ∈ r.tup do {
    if t[a] = e then {
      q.tup = q.tup ∪ {t};
    }
  }
  return q;
}
```

Implementation in RISCAL

Definition

Let $\mathcal{A} = \{A_1, \dots, A_n\}$ be an attribute schema. The operation $\pi : \text{Relation}_{\mathcal{A}} \times \mathcal{P}(\mathcal{A}) \rightarrow \bigcup_{\mathcal{A}' \subseteq \mathcal{A}} \text{Relation}_{\mathcal{A}'}$ defined by

$$\pi(r, \mathcal{B}) = \{t|_{\mathcal{B}} : t \in r\}$$

is called projection of r on \mathcal{B} .

RISCAL

```
pred project_spec(s:Relation, r:Relation, columns:Array[N,Length]) ⇔
s.len = |{i | i:Attribute with columns[i] ≠ N}| ∧ (∀ tr:Row. tr ∈ r.tup ⇒
  ∃ ts:Row. ts ∈ s.tup ∧ ∀ i:Attribute. i < s.len ⇒ ts[i]=tr[columns[i]]);
```

Implementation in RISCAL

```
fun project1(r:Relation, columns:Array[N,Length]):Relation
requires ( $\exists$  i:Attribute.  $\forall$  j:Attribute. (j>i  $\Rightarrow$  columns[j] = N)  $\wedge$  (j $\leq$ i  $\Rightarrow$  columns[j] < r.
  len));
= choose s:Relation with project_spec(s,r,columns);
```

```
proc project2(r:Relation, columns:Array[N,Length]):Relation
requires ( $\exists$  i:Attribute.  $\forall$  j:Attribute. (j>i  $\Rightarrow$  columns[j] = N)  $\wedge$  (j $\leq$ i  $\Rightarrow$  columns[j] < r.
  len));
ensures project_spec(result,r,columns); {
```

```
  var l:Length := |{i | i:Attribute with columns[i]  $\neq$  N}|;
  var q:Relation :=  $\langle$ len: l, tup: choose s:Set[Row] with |s|=0);
```

```
  for t  $\in$  r.tup do {
    var tn:Row := Array[N,Element](0);

    var j:Length := 0;
    for var i:Length := 0; i<N; i:=i+1 do {
      if columns[i]  $\neq$  N then {
        tn[j] := t[columns[i]];
        j := j+1;
      }
    }
    q.tup := q.tup  $\cup$  {tn};
  }
```

```
  return q;
```

```
}
```

Model Checking Algorithms and Theorems

- For each function/procedure:
 - execute the function and
 - check the postcondition for every valid input
- For each theorem: Check validity

Both tasks can be accomplished by the RISCAL task "Execute operation".

Model Checking Algorithms and Theorems

Theorem 1

Let \mathcal{A} be an attribute schema, $A \in \mathcal{A}$, $a \in \text{dom}(A)$ and $r, s \in \text{Relation}_{\mathcal{A}}$ relations. Then

$$\sigma(r\gamma s, A, a) = \sigma(r, A, a)\gamma\sigma(s, A, a)$$

holds for all $\gamma \in \{\cup, \cap, \setminus\}$.

```
theorem select_union_equiv(r1:Relation, r2:Relation, a:Attribute, e:Element)
requires a < r1.len ^ a < r2.len ^ union_compatible(r1,r2) ^ |r1.tupl| + |r2.tupl| ≤ M; ⇔
select2(rUnion(r1,r2),a,e) = rUnion(select2(r1, a, e), select2(r2, a, e));
```

```
theorem select_intersect_equiv(r1:Relation, r2:Relation, a:Attribute, e:Element)
requires a < r1.len ^ a < r2.len ^ union_compatible(r1,r2); ⇔
select2(rIntersect(r1,r2),a,e) = rIntersect(select2(r1, a, e), select2(r2, a, e));
```

```
theorem select_minus_equiv(r1:Relation, r2:Relation, a:Attribute, e:Element)
requires a < r1.len ^ a < r2.len ^ union_compatible(r1,r2); ⇔
select2(rMinus(r1,r2),a,e) = rMinus(select2(r1, a, e), select2(r2, a, e));
```

Theorem 2

Let \mathcal{A} be an attribute schema, $A, B \in \mathcal{A}$, $a \in \text{dom}(A)$, $b \in \text{dom}(B)$ and $r \in \text{Relation}_{\mathcal{A}}$ a relation. Then

$$\sigma(\sigma(r, B, b), A, a) = \sigma(r, A, a) \cap \sigma(r, B, b)$$

holds.

```
theorem select_intersect_comp(r:Relation, a:Attribute, e:Element, b:Attribute, f:Element)
requires a < r.len ^ b < r.len; ⇔
select2(select2(r, a, e), b, f) = rIntersect(select2(r, a, e), select2(r, b, f));
```

Theorem 3

Let $\mathcal{A} = \{A_1, \dots, A_n\}$ be an attribute schema, $A, B \in \mathcal{A}$ and $r \in \text{Relation}_{\mathcal{A}}$ a relation. Then

$$\pi(r, \{A, B\}) \subseteq C(\pi(r, \{A\}), \pi(r, \{B\}))$$

holds.

```
theorem project_cartesian_subset(r:Relation, a:Attribute, b:Attribute)
requires a < r.len ^ b < r.len ^ 2*|r.tup| ≤ M; ⇔
project2(r, attributes(List!node(a, List!node(b, List!nil)))) .tup
  ⊆ cartesian2(project2(r, attributes(List!node(a, List!nil))),
               project2(r, attributes(List!node(b, List!nil)))) .tup;
```

Theorem 3

Let \mathcal{A}, \mathcal{B} be attribute schemas, $A \in \mathcal{A}$, $B \in \mathcal{B}$ and $r \in \text{Relation}_{\mathcal{A}}$, $s \in \text{Relation}_{\mathcal{A}}$ relations. Then

$$\bowtie(r, s, A, B) \subseteq C(r, s)$$

holds.

```
theorem join_cartesian_subset(r1:Relation, r2:Relation, n1:Attribute, n2:Attribute)
requires n1 < r1.len ^ n2 < r2.len ^ r1.len + r2.len ≤ N ^ |r1.tup| * |r2.tup| ≤ M; ⇔
join3(r1, r2, n1, n2).tup ⊆ cartesian3(r1, r2).tup;
```

An SQL-like wrapper language

Grammar

```
 $L \in List$   
 $Q \in Query$   
 $L ::= () \mid (A, L)$   
 $Q ::= from(T)$   
      |  $on(A_1, A_2, Q_1, Q_2)$   
      |  $where(A, E, Q)$   
      |  $select(L, Q)$   
      |  $add(Q_1, Q_2)$   
      |  $inters(Q_1, Q_2)$   
      |  $minus(Q_1, Q_2)$   
      |  $cart(Q_1, Q_2)$ 
```

- The natural variables T and A identify tables and attributes, whereas E denotes either 0 or 1.
- The language SQL_{min} consists of all expressions in the syntactical domain $Query$.

An SQL-like wrapper language

```
rectype(L) List = nil | node(Attribute, List);
rectype(D) Query = from(TableId)
  | on(Attribute, Attribute, Query, Query)
  | where(Attribute, Element, Query)
  | select(List, Query)
  | add(Query, Query)
  | inters(Query, Query)
  | minus(Query, Query)
  | cart(Query, Query);
```

An SQL-like wrapper language

Definition

Let $db : \text{TableId} \rightarrow \bigcup_{i=1}^n \text{Relation}_{\mathcal{A}_i}$ be a database. We define the denotational semantics $\llbracket \cdot \rrbracket_{db} : \text{Query} \rightarrow \text{Relation}_{\mathcal{A}}$ of SQL_{\min} by recursive induction in the following way:

$$\llbracket \text{from}(T) \rrbracket_{db} := db(T)$$

$$\llbracket \text{on}(A_1, A_2, Q_1, Q_2) \rrbracket_{db} := \bowtie (\llbracket Q_1 \rrbracket_{db}, \llbracket Q_2 \rrbracket_{db}, A_1, A_2)$$

$$\llbracket \text{where}(A, E, Q) \rrbracket_{db} := \sigma(\llbracket Q \rrbracket_{db}, A, E)$$

$$\llbracket \text{select}(L, Q) \rrbracket_{db} := \pi(\llbracket Q \rrbracket_{db}, L)$$

$$\llbracket \text{add}(Q_1, Q_2) \rrbracket_{db} := \llbracket Q_1 \rrbracket_{db} \cup \llbracket Q_2 \rrbracket_{db}$$

$$\llbracket \text{inters}(Q_1, Q_2) \rrbracket_{db} := \llbracket Q_1 \rrbracket_{db} \cap \llbracket Q_2 \rrbracket_{db}$$

$$\llbracket \text{minus}(Q_1, Q_2) \rrbracket_{db} := \llbracket Q_1 \rrbracket_{db} \setminus \llbracket Q_2 \rrbracket_{db}$$

$$\llbracket \text{cart}(Q_1, Q_2) \rrbracket_{db} := C(\llbracket Q_1 \rrbracket_{db}, \llbracket Q_2 \rrbracket_{db})$$

An SQL-like wrapper language

```
proc query(db:Database, q:Query):Relation {
  var r:Relation := match q with {
    from(tid:TableId) -> db[tid];
    on(n1:Attribute, n2:Attribute, q1:Query, q2:Query) -> join(query(db, q1), query(db,
      q2), n1, n2);
    where(a:Attribute, e:Element, q:Query) -> select(query(db, q), a, e);
    select(a:List, q:Query) -> project(query(db, q), attributes(a));
    add(q1:Query, q2:Query) -> rUnion(query(db, q1), query(db, q2));
    inters(q1:Query, q2:Query) -> rIntersect(query(db, q1), query(db, q2));
    minus(q1:Query, q2:Query) -> rMinus(query(db, q1), query(db, q2));
    cart(q1:Query, q2:Query) -> cartesian(query(db, q1), query(db, q2));
  };
  return r;
}
```

Translation into standardized SQL

A SQL query is always executed in the same order, i.e. we can sort the clauses in the way

$$\underset{4}{\text{FROM}} \succ \underset{3}{\text{ON}} \succ \underset{2}{\text{WHERE}} \succ \underset{1}{\text{SELECT}},$$

where $C_1 \succ C_2$ means that clause C_1 is evaluated before clause C_2 . If the corresponding operations of SQL_{\min} are composed in this order, they are equivalent to a single SQL query without subqueries.

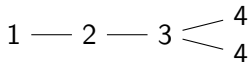
Translation into standardized SQL

We can list every possibility, because

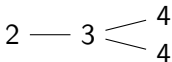
- the leaves of any AST must contain FROM clauses and
- under the assumption that the tree nodes of an AST contain clauses of strictly ascending priority, there are maximal 4 layers.

This equates to the eighth "atomic" SQL expressions.

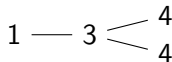
Translation into standardized SQL



```
SELECT <attributes>  
FROM <table>  
JOIN <table>  
ON <attribute> = <attribute>  
WHERE <attribute> = <value>;
```



```
SELECT *  
FROM <table>  
JOIN <table>  
ON <attribute> = <attribute>  
WHERE <attribute> = <value>;
```



```
SELECT <attributes>  
FROM <table>  
JOIN <table>  
ON <attribute>  
= <attribute>;
```



```
SELECT <attributes>  
FROM <table>  
WHERE <attribute> = <value>;
```



```
SELECT *  
FROM <table>  
JOIN <table>  
ON <attribute> = <attribute>;
```



```
SELECT *  
FROM <table>  
WHERE <attribute>  
= <value>;
```



```
SELECT <attributes>  
FROM <table>;
```



```
SELECT *  
FROM <table>;
```

Translation into standardized SQL

```
proc query_to_sql(q0:Query):()  
proc translate_on(n1:Attribute, n2:Attribute, q0_1:Query, q0_2:Query):()  
proc translate_cart(q0_1:Query, q0_2:Query):()  
proc translate_where(a:Attribute, e:Element, q0:Query):()  
proc translate_from(tid:TableId):()  
proc translate_blank(q0:Query):()  
proc translate_select(l>List, q0:Query):()
```

Translation into standardized SQL

```
proc query_to_sql(q0:Query):() {
  match q0 with {
    add(q1_1:Query, q1_2:Query) -> {
      translate_add(q1_1, q1_2);
    }
    inters(q1_1:Query, q1_2:Query) -> {
      translate_inters(q1_1, q1_2);
    }
    minus(q1_1:Query, q1_2:Query) -> {
      print "Not implemented";
    }
    cart(q1_1:Query, q1_2:Query) -> {
      translate_cart(q1_1, q1_2);
    }
    select(l>List, q1:Query) -> {
      translate_select(l, q1);
    }
    where(a:Attribute, e:Element, q1:Query) -> {
      print "SELECT *";
      translate_where(a, e, q1);
    }
    on(n1:Attribute, n2:Attribute, q1_1:Query, q1_2:Query) -> {
      print "SELECT *";
      translate_on(n1, n2, q1_1, q1_2);
    }
    from(tid:TableId) -> {
      print "SELECT *";
      translate_from(tid);
    }
  }
}
```

Translation into standardized SQL

RISCAL

```
Query!on(0,0,  
  Query!select(  
    List!node(0,  
      List!node(2, List!nil)),  
    Query!where(1,1,  
      Query!from(0)  
    )  
  ),  
  Query!from(1)  
);
```

query_to_sql

```
SELECT * FROM (  
  SELECT 0, 2  
  FROM 0  
  WHERE 1 = 1  
)  
JOIN 1  
ON 0 = 0
```

Sqlite

```
SELECT * FROM (  
  SELECT "0"."0", "0"."2"  
  FROM "0"  
  WHERE "0"."1" = 1  
) AS "_0"  
JOIN "1"  
ON "_0"."0" = "1"."0";
```


Current progress and further goals

Done:

- Formalization and verification with a Model Checker
- Design of a Language

Work in progress:

- Theorem proving with RISCTP
- Various formal improvements