

A Crash Course in Database Queries

Jan Van den Bussche
Hasselt University

joint work with Stijn Vansummeren, Dirk Van Gucht

Query languages are programming languages!

[Atkinson–Buneman–Cardelli–Maier–Ohuri–Sheard–Stemple
–Stonebraker–Tannen]

⇒ Database queries. . .

- are programs
- can crash
- can be ill-typed
- are polymorphic
- can be treated as data (metadata, reflection)

Motivation

Lowell self-assessment:

“We recommend that database researchers increase their focus on the integration of text, data, code, and streams.”

Asilomar Report:

“Ever more complex application environments have increased the need to integrate programs and data.”

What are the basic theoretical questions concerning flexible operation of queries “out of the box”? Of treating queries as data?

How to apply/adapt programming-language ideas to query languages?

Query languages

Relational algebra

Nested relational calculus

XQuery

Relational algebra: syntax

$$\begin{array}{l} e ::= x \quad \text{(relation variable)} \\ \quad | r \quad \text{(constant relation)} \\ \quad | e \cup e \\ \quad | e - e \\ \quad | e \times e \\ \quad | \sigma_{A=B}(e) \\ \quad | \pi_{A,\dots,B}(e) \\ \quad | \rho_{A/B}(e) \end{array}$$

Example expression:

$$x - (\pi_A(x) \times \rho_{C/B}(y))$$

Heterogeneous relations

A relation is a finite set of tuples

A tuple is a mapping t from some relation scheme to \mathbb{V}

- \mathbb{V} : universe of atomic values
- relation scheme: finite set of attributes
- call relation scheme of t , the type of t

$(A: 1, B: 2)$
$(A: 3, C: 4)$
$(D: 4)$

Relational algebra: semantics

Expression $e(x, \dots, y)$ can be applied to relations r, \dots, s

Operational semantics: rewrite ground expression $e' = e(r, \dots, s)$
until you end up with a relation r'

Rewriting $e' \rightarrow r'$ defined by inference rules

Inference rule for union

if $e_1 \rightarrow r_1$
and $e_2 \rightarrow r_2$
then $e_1 \cup e_2 \rightarrow r_1 \cup r_2$

Written like fraction:

$$\frac{e_1 \rightarrow r_1 \quad e_2 \rightarrow r_2}{e_1 \cup e_2 \rightarrow r_1 \cup r_2}$$

Difference:

$$\frac{e_1 \rightarrow r_1 \quad e_2 \rightarrow r_2}{e_1 - e_2 \rightarrow r_1 \setminus r_2}$$

Cartesian product

$$\frac{\begin{array}{l} e_1 \rightarrow r_1 \\ e_2 \rightarrow r_2 \\ \forall t_1 \in r_1 : \forall t_2 \in r_2 : \text{type}(t_1) \cap \text{type}(t_2) = \emptyset \end{array}}{e_1 \times e_2 \rightarrow \{t_1 \cup t_2 \mid t_1 \in r_1 \ \& \ t_2 \in r_2\}}$$

Selection, projection, renaming

$$\frac{e \rightarrow r' \quad \forall t \in r' : A, B \in \text{type}(t)}{\sigma_{A=B}(e) \rightarrow \{t \in r' \mid t(A) = t(B)\}}$$

$$\frac{e \rightarrow r' \quad \forall t \in r' : A, \dots, B \in \text{type}(t)}{\pi_{A, \dots, B}(e) \rightarrow \{\pi_{A, \dots, B}(t) \mid t \in r'\}}$$

$$\frac{e \rightarrow r' \quad \forall t \in r' : A \in \text{type}(t) \ \& \ B \notin \text{type}(t)}{\rho_{A/B}(e) \rightarrow \{\rho_{A/B}(t) \mid t \in r'\}}$$

Example derivation

Evaluate $x - (\pi_A(x) \times \rho_{C/B}(y))$

$$\begin{array}{r}
 r \rightarrow \begin{array}{|l} (A: 1, B: 2) \\ (A: 1, B: 3) \\ (A: 2, B: 4) \\ (A: 2, B: 5) \end{array} \\
 \hline
 \pi_A(r) \rightarrow \begin{array}{|l} (A: 1) \\ (A: 2) \end{array} \\
 \hline
 \end{array}
 \qquad
 \begin{array}{r}
 s \rightarrow \begin{array}{|l} (C: 2) \\ (C: 3) \\ (C: 4) \end{array} \\
 \hline
 \rho_{C/B}(s) \rightarrow \begin{array}{|l} (B: 2) \\ (B: 3) \\ (B: 4) \end{array} \\
 \hline
 \end{array}$$

$$\begin{array}{r}
 r \rightarrow r \qquad \pi_A(r) \times \rho_{C/B}(s) \rightarrow \begin{array}{|l} (A: 1, B: 2) \\ (A: 1, B: 3) \\ (A: 1, B: 4) \\ (A: 2, B: 2) \\ (A: 2, B: 3) \\ (A: 2, B: 4) \end{array} \\
 \hline
 r - (\pi_A(r) \times \rho_{C/B}(s)) \rightarrow \begin{array}{|l} (A: 2, B: 5) \end{array}
 \end{array}$$

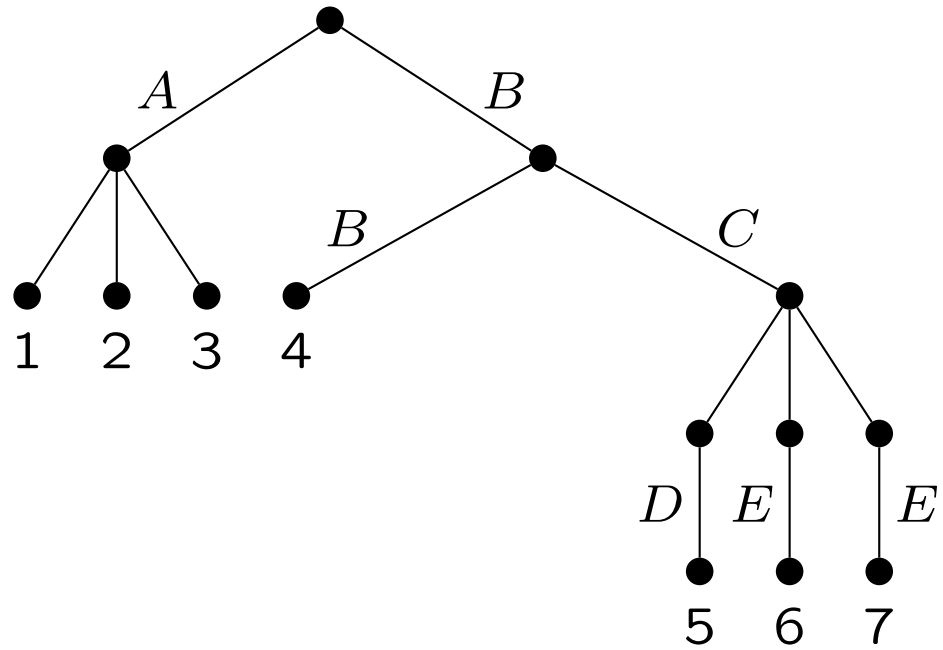
Nested relations (complex objects)

A complex object is:

- an atomic value
- a tuple of complex objects
- a finite set of complex objects

$$(A: \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 3 \\ \hline \end{array}, B: (B: 4, C: \begin{array}{|c|} \hline (D: 5) \\ \hline (E: 6) \\ \hline (E: 7) \\ \hline \end{array}))$$

Complex object tree



$$(A: \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 3 \\ \hline \end{array}, B: (B: 4, C: \begin{array}{|c|} \hline (D: 5) \\ \hline (E: 6) \\ \hline (E: 7) \\ \hline \end{array}))$$

Nested relational calculus

[BNTW]

The relational algebra of complex objects

Syntax:

$e ::=$	x	(variable)
	r	(constant object)
	\emptyset	
	$\{e\}$	
	$e \cup e$	
	$\bigcup e$	
	$e.A$	
	$(A: e, \dots, B: e)$	
	for $z \in e$ return e	
	if $e \text{ eq } e$ then e else e	
	if $e = e$ then e else e	

Example NRC expression

∪ for $u \in x$ return
 ∪ for $v \in y$ return
 if $u.B$ eq $v.B$
 then $\{(A: u.A, B: u.B, C: v.C)\}$
 else \emptyset

Operational semantics of for-loop

$$\frac{e_1 \rightarrow r \quad r \text{ is a set} \quad \forall t \in r : e_2(t) \rightarrow s_t}{\text{for } z \in e_1 \text{ return } e_2(z) \rightarrow \{s_t \mid t \in r\}}$$

Equality test

$$\frac{e_1 \rightarrow r_1 \quad e_2 \rightarrow r_2 \quad e_3 \rightarrow r_3 \quad r_1 = r_2}{\text{if } e_1 = e_2 \text{ then } e_3 \text{ else } e_4 \rightarrow r_3}$$

$$\frac{e_1 \rightarrow r_1 \quad e_2 \rightarrow r_2 \quad e_4 \rightarrow r_4 \quad r_1 \neq r_2}{\text{if } e_1 = e_2 \text{ then } e_3 \text{ else } e_4 \rightarrow r_4}$$

$e_1 \text{ eq } e_2$: similar, but r_1 and r_2 must be atomic values

Positive NRC: has only $e_1 \text{ eq } e_2$

Queries can crash

$\pi_A(x)$ crashes on $r = \boxed{\begin{array}{l} (A: 1, B: 2) \\ (B: 3) \end{array}}$

$x \times y$ crashes on $r = \boxed{\begin{array}{l} (A: 1) \\ (B: 2) \end{array}}$ and $s = \boxed{(B: 3, C: 4)}$

Similar for NRC

“Well-defined” = “does not crash”

Cannot expect an expression to be well-defined on all inputs

\Rightarrow consider types

The well-definedness problem

(Relational algebra version)

Recall: type = relation scheme = finite set of attributes

A relation r has type R if all its tuples do

– notation $r : R$

The well-definedness problem:

Input: expression $e(x, \dots, y)$ and types R, \dots, S

Decide: is $e(r, \dots, s)$ well-defined on all inputs $r : R, \dots, s : S$?

Some immediate observations

Well-definedness is undecidable for relational algebra:

- take $R = \{A, B\}$

- for well-defined e :

$$e(x) \text{ satisfiable} \quad \Leftrightarrow \quad \pi_A \pi_B(x \times \pi_\emptyset(e)) \text{ ill-defined}$$

Decidable for positive relational algebra (without difference):

- just keep track of possible types that can occur

- monotonicity

- can add $\sigma_{A \neq B}$

The NRC case

Types:

- atom
- $(A: R, \dots, B: S)$
- $\{T\}$

Well-definedness still undecidable for full NRC

Well-definedness for positive NRC

Still decidable

- small model property for ill-definedness
- monotonicity
- coNEXPTIME-hard (satisfiability, [Koch])

Singleton extraction

$$\frac{e \rightarrow \{t\}}{\text{extract}(e) \rightarrow t}$$

OQL, XQuery, but also SQL:

```
select ..., (Q), ... from ... where ...
```

or

```
select ... from ... where A = (Q)
```

crashes unless subquery returns a singleton (“scalar subquery”)

Well-definedness for positive NRC + extract

Is undecidable

- $\text{extract}(\{e_1, e_2\})$ well-defined iff e_1, e_2 equivalent
- equivalence of positive NRC is undecidable!

(Satisfiability still decidable.)

Decidable for lists, bags (see also XQuery)

XQuery (w/o recursion)

Value = list of atoms and tree nodes

- underlying store of XML trees (node-labeled by atoms)

Language:

$e ::= x$	(variable)
a	(constant atom)
$()$	
for $z \in e$ return e	
if e then e else e	
let $z := e$ in e	
$f(e, \dots, e)$	(operators)

XQuery operators

Operators: Element construction, list functions, axes, tests, . . .

– can crash!

e.g. $\text{element}\{e_1\}\{e_2\} \Rightarrow e_1$ must be singleton

XQuery well-definedness

Types: bounded-depth regular expression types

- XML Schema (extended DTD) w/o recursion

For well-behaved (generic, monotonic, local) operators, XQuery well-definedness is decidable

Caveat: automatic coercions

- atomization: tree \rightarrow atom
- difference between NRC(=) and NRC(eq) is blurred
- undecidability

Semantic type-checking

Input: expression $e(x, \dots, y)$, well-defined under types R, \dots, S ;
additional type T

Decide: is $e(r, \dots, s)$ always of type T , for all inputs
 $r : R, \dots, s : S$?

- RA, NRC: same story as well-definedness
- NRC + extract: still decidable!
- XQuery: undecidable!

Static type checking

Since input types are known, can try to derive statically:

- can expression crash?
- if not, what is output type?

⇒ Inference rules that derive $\Gamma \vdash e : T$

– Γ : the given input types

- relational algebra
- NRC
- XQuery [XQuery formal semantics; Ghelli et al., JFP 2006]

Static type system of relational algebra

$$\frac{\Gamma \vdash e_1 : T \quad \Gamma \vdash e_2 : T}{\Gamma \vdash e_1 \cup e_2 : T}$$

$$\frac{\Gamma \vdash e_1 : T \quad \Gamma \vdash e_2 : T}{\Gamma \vdash e_1 - e_2 : T}$$

$$\frac{\Gamma \vdash e_1 : R \quad \Gamma \vdash e_2 : S \quad R \cap S = \emptyset}{\Gamma \vdash e_1 \times e_2 : R \cup S}$$

$$\frac{\Gamma \vdash e : T \quad A, B \in T}{\Gamma \vdash \sigma_{A=B}(e) : T}$$

$$\frac{\Gamma \vdash e : T \quad A, \dots, B \in T}{\Gamma \vdash \pi_{A, \dots, B}(e) : \{A, \dots, B\}}$$

$$\frac{\Gamma \vdash e : T \quad A \in T \quad B \notin T}{\Gamma \vdash \rho_{A/B}(e) : (R \setminus \{A\}) \cup \{B\}}$$

Soundness and completeness of type systems

If e is well-typed ($e : T$ can be derived) then e is well-defined, and output type is always T

Not vice versa: type checking is only sound, not complete

Nevertheless we have expressive completeness:

- a well-defined expression with known output type can always be equivalently rewritten into a well-typed expression

[Vansummeren]

Expressive completeness of static type checking

- given types R, \dots, S
- given type T
- given expression $e(x, \dots, y)$
 - e well-defined under R, \dots, S
 - e 's output type is always T

Then there always exists e' equivalent to e on all inputs $r : R, \dots, s : S$, and

$$\boxed{x : R, \dots, y : S \vdash e' : T}$$

Holds for RA, NRC — for XQuery?

Simple example of expressive completeness

Consider $x : \{A, B\}$ and $y : \{A, C\}$

Then

$$\pi_A(x \cup y)$$

is equivalent to

$$\pi_A(x) \cup \pi_A(y)$$

Less trivial for queries involving difference

- Encode, in RA/NRC, heterogeneous relations/objects by homogeneous ones.

Note

Expressive completeness is immediate for well-typed languages that are Turing-complete

- Java, Haskell, . . .
- Simply-typed lambda calculus: not expressive complete

Polymorphism

$\hat{\pi}_A$: project out attribute A

\bowtie : natural join

Typing rules:

$$\frac{e : R \quad A \in R}{\hat{\pi}_A(e) : R \setminus \{A\}}$$

$$\frac{e_1 : R_1 \quad e_2 : R_2}{e_1 \bowtie e_2 : R_1 \cup R_2}$$

For any given input types, $\hat{\pi}$ and \bowtie can be expressed using the other RA operators

But not by one expression that works for all possible input types!

The ultimate polymorphic, typed, query language?

Type inference

Input: expression e

Output: all Γ and T for which $\Gamma \vdash e : T$

Output is usually infinite; need some simple kind of finite representation

Output can be empty as well (untypeable expression, e.g., $\pi_A \pi_B(x)$)

\Rightarrow type formulas

Polytypes with kinding [Ohori–Buneman]

Polytype = type with type variables (ML, Hindley–Milner)

NRC:

for $u \in x$ return
 for $v \in y$ return
 $(C : u.A, D : v.B)$

Type formula with type variables and kinding

$$\begin{array}{l} x : \{\alpha\} \\ y : \{\beta\} \\ \text{kind}(\alpha) = (A : \gamma) \mapsto (C : \gamma, D : \delta) \\ \text{kind}(\beta) = (B : \delta) \end{array}$$

Insufficient to represent type inference of $\hat{\pi}$ or ρ

Row variables

[Rémy]

$$\hat{\pi}_A(x) \cup y$$

Type formula with row variables and forbidden attributes:

$$\begin{array}{l} x : \{A\} \cup \alpha \\ y : \alpha \end{array} \mapsto \rho$$
$$\text{forbidden}(\alpha) = \{A\}$$

Insufficient to represent type inference of \times

Type inference for full relational algebra

1. Use multiple row variables, stand for disjoint types
2. Generalize required, forbidden attributes to boolean constraints

$$e = \sigma_{B=C}(\rho_{A/B}(x) \times y)$$

$$\begin{array}{l} x : \alpha \\ y : \beta \\ A : x \\ B : \neg x \wedge \neg y \\ C : x \vee y \end{array} \quad \mapsto \quad \begin{array}{l} e : \alpha \cup \beta \\ A : y \\ B : \text{true} \\ C : \text{true} \end{array}$$

Type formulas can become exponentially long, but typeability is in NP (complete)

General approach to type inference

1. Universe of all types, with appropriate constraints and operations on types, forms a logical model \mathcal{M}
2. Formulate quantifier-free formula over type variables, stating constraints on input types of e , for e to be well-typed
 - Type formula!
3. Existential theory of \mathcal{M} is decidable (typeability)

Can play this game for full NRC with $\hat{\pi}$, ρ , \boxtimes , \times , \dots

Complexity for NRC does not become worse than for relational algebra (NP-complete)

Type reflection

java.lang.reflection package:

- inspect type (class) of an object
- type information becomes data!

Schema querying ('90s)

Semistructured data, XML: distinction between type information and data is blurred

- can validate XML Schema using (recursive) XQuery

Downside: is XPath expression $x/a=5$ false because x has no a -child? Or because x has an a -child, but it is not 5?

Example: transposition of a relation [GL,WR]

<i>A</i>	<i>B</i>	<i>C</i>
1	<i>d</i>	<i>e</i>
2	<i>f</i>	<i>g</i>

 ↔

<i>A</i>	1	2
<i>B</i>	<i>d</i>	<i>f</i>
<i>C</i>	<i>e</i>	<i>g</i>

⇒ Attributes as data values

Expressions as data

E.g., workload log:

user	datetime	query
john	20070612T1030	select ... from ... where ...
mary	20070611T2316	...
	:	

Integration of program logic and data

- System catalog: VIEWS table
- QUEL as a data type
- Oracle EXPRESSION type
- workload monitoring
- publish-subscribe
- workflow management
- software engineering

Querying a database containing queries

- Hotspots: which subqueries are often used?
 - syntactic
- Semantic: which queries return no answer?
 - semantic
- View maintenance: how do the query answers change under this update?
 - syntactic & semantic

Two approaches: (i) decomposition; (ii) ADT

Decomposition-based approach

Use standard query language for syntactic manipulations

⇒ Stored expressions cannot be represented as atomic values

- must be decomposed

Many ways to do this:

- XML: syntax tree
 - XQueryX
- Relational: decomposition

MetaSQL

SQL/XML

- queries are stored in XML columns

Add EVAL function to SQL

- analogy with Lisp, Scheme

E.g. workloads: Log(user, datetime, query)

```
select query, A, B
from Log
where (A,B) in EVAL(rewrite(query,update)) MINUS EVAL(query)
```

Similarly can add EVAL to RA, or XQuery

Does EVAL add power?

Data complexity of EVAL = evaluation complexity of relational algebra evaluation:

Input: database D and expression e

Output: $e(D)$

PSPACE-complete \gg LOGSPACE (plain relational algebra)

But what about standard generic queries?

- expressions not in input
- dynamic generation and evaluation of expressions only as auxiliary querying tool

Transitive closure in XQuery + EVAL

Table $R(A, B)$ in XML: list D of $R(A, B)$ -elements

$TC(D) = EVAL(\underline{\text{construct}}(D)):$

$\text{construct}(D) = E_1, \dots, E_n$ with $n = \text{count}(D//T)$

with $E_j:$

for t_1 in $D//T$, ..., t_j in $D//T$ return
if every z in $((t_1/B=t_2/A), \dots, (t_{j-1}/B=t_j/A))$
satisfies $z=fn:true()$ then
 $\text{element}(T)\{t_1/A, t_j/B\}$ else $()$

In relational algebra, on databases w/o stored expressions,
EVAL = for-loops

Type-safe reflection

EVAL: fragile, can crash easily

Idea: two-level type system [MetaML], e.g.:

$(A: \text{atom}, B: \text{atom}, C: \langle\{(D: \text{atom}, E: \text{atom})\}\rangle)$

\Rightarrow EVAL can be typed

ADT approach: also provide typed repertoire of syntactic manipulation operators

e.g. substitute all occurrences of relation variable $x : T$ by the expression $e : T$

Less powerful \Rightarrow add polymorphism?

Conclusions

Flexible operation of queries “out of the box”

- Well-definedness (better algorithms?)
- Expressive completeness of type systems
- Polymorphism (design of query languages?)
- Meta-querying: querying queries (design of languages?)

Integration of programs (queries) and data