

A Uniform Calculus for Collection Types

Leonidas Fegaras
Oregon Graduate Institute

Problems with ODMG-93 OQL

- Implicit type conversions:

```
select  $e$   
from  $x_1$  in  $e_1, \dots, x_n$  in  $e_n$   
where  $pred$ 
```

OQL converts each e_i into a set before the join!
 $flatten(x)$ is even messier.

- Missing operations:

```
select list  $e.B$   
from  $e$  in (sort  $s$  in  $x$  by  $s.A$ )
```

- Redundancy:

$flatten(x) = \text{select } a \text{ from } s \text{ in } x, a \text{ in } s$

Fundamental questions

- Is the join of a list with a set meaningful?
- If so, what is the result type of this join?
- More generally, what are the precise semantics of queries over multiple collection types?
- How do we optimize such queries?

Extending the (nested) relational algebra is not necessarily a good solution

- some operations may have many instances and interpretations:
e.g. join a list with a bag and return a set,
e.g. $\text{flatten}(x)$ returns a set if x is a set of sets,
a bag if x is a bag of sets, etc;
(EREQ/AQUA has about 100 operators!)
- the more bulk operations you support the harder the optimization task becomes;
(need $\mathcal{O}(n^2)$ rules for n operations)
- testing program equivalences is more feasible for a small number of primitives;
- adding a new collection type will affect many operators.

Claims

- a query algebra doesn't have to resemble the query language or the physical algorithms;
- a query algebra should consist of a small number of powerful, generic, non-overlapping, operators (just one if possible).

But:

How these powerful operators are mapped into physical algorithms?

Motivation: Example from OQL

Cities : $set(\langle name : string, hotels : bag(\dots), places_to_visit : list(\dots) \rangle)$

OQL query:

```
select distinct h.name
from c in Cities,
      h in c.hotels,
      p in c.places_to_visit,
where c.name = "Madison" and h.name = p.name
```

In NF²-style algebraic form:

```
project(flatten(map(filter(Cities,  $\lambda c. c.name = \text{"Madison"}$ ),
                   $\lambda c. \mathbf{join}(c.hotels, c.places\_to\_visit,$ 
                               $\lambda(h, p). h.name = p.name)$ )),
         $\lambda(h, p). h.name$ )
```

In my calculus (comprehension syntax):

```
 $set\{ h.name \mid c \leftarrow \text{Cities}, h \leftarrow c.hotels, p \leftarrow c.places\_to\_visit,$ 
         $c.name = \text{"Madison"}, h.name = p.name \}$ 
```

Another example:

```
select h.name
from h in (select c.hotels
               from c in Cities
               where c.name = "Madison")
where exists r in h.rooms : (r.bed# = 3)
```

```
bag{ h.name | h ← bag{ c.hotels | c ← Cities, c.name = "Madison" },
     some{ r.bed# = 3 | r ← h.rooms } }
```

Monoids

$$(T, \text{zero}[T], \text{merge}[T])$$

is a **monoid** if $\text{merge}[T]$ is **associative** with a zero element:

$$\begin{aligned} \text{merge}[T](\text{merge}[T](x, y), z) &= \text{merge}[T](x, \text{merge}[T](y, z)) \\ \text{merge}[T](\text{zero}[T], x) &= \text{merge}[T](x, \text{zero}[T]) = x \end{aligned}$$

Optional properties:

- C) **Commutativity:** $\text{merge}[T](x, y) = \text{merge}[T](y, x)$
- I) **Idempotence:** $\text{merge}[T](x, x) = x$

My approach: Every monoid is **predefined!**

<i>free monoid T</i>	$\text{zero}[T]$	$\text{unit}[T](a)$	$\text{merge}[T]$	C/I
list	$[\]$	$[a]$	append	
set	$\{\}$	$\{a\}$	\cup	CI
bag	$\{\{\}$	$\{\{a\}$	\uplus	C
ordered-set	$\{\}\}$	$\{a\}$	\cup	I

<i>simple monoid T</i>	type	$\text{zero}[T]$	$\text{unit}[T](a)$	$\text{merge}[T]$	C/I
sum	int	0	a	+	C
prod	int	1	a	*	C
max	int	0	a	max	CI
some	boolean	false	a	\vee	CI
all	boolean	true	a	\wedge	CI

Monoid Type:

$type = T$ (T is a simple monoid)

$type = T(type')$ (T is a free monoid)

$type = \langle a_1 : type_1, \dots, a_n : type_n \rangle$ (aggregation)

e.g.

Cities : $set(\langle name : string, hotels : bag(\dots), places_to_visit : list(\dots) \rangle)$

A monoid type is a monoid.

e.g.

$$\begin{aligned} \text{zero}[T(type)] &= \text{zero}[T] \\ \text{unit}[T(type)](x) &= \text{unit}[T](\text{unit}[type](x)) \\ \text{merge}[T(type)](x, y) &= \text{merge}[T](x, y) \end{aligned}$$

Monoid Homomorphisms

$\text{hom}[T, S](f)$ is a **homomorphism** from monoid T to monoid S :

$$\begin{aligned}\text{hom}[T, S](f) (\text{zero}[T]) &= \text{zero}[S] \\ \text{hom}[T, S](f) (\text{unit}[T](a)) &= f(a) \\ \text{hom}[T, S](f) (\text{merge}[T](x, y)) &= \begin{cases} \text{merge}[S](\text{hom}[T, S](f)(x), \\ \text{hom}[T, S](f)(y)) \end{cases}\end{aligned}$$

If $T = \text{list}$ and $S = \text{set}$ then:

$$\begin{aligned}\text{hom}[\text{list}, \text{set}](f) ([]) &= \{\} \\ \text{hom}[\text{list}, \text{set}](f) ([a]) &= f(a) \\ \text{hom}[\text{list}, \text{set}](f) (\text{append}(x, y)) \\ &= \text{hom}[\text{list}, \text{set}](f)(x) \cup \text{hom}[\text{list}, \text{set}](f)(y)\end{aligned}$$

e.g. if $x = [a_0, a_1, \dots, a_n]$ then:

$$\text{hom}[\text{list}, \text{set}](f)(x) = f(a_0) \cup f(a_1) \cup \dots \cup f(a_n)$$

If $S = T$ then $\text{hom}[T, T](f)$ is $\text{flattenmap}(f)$ in [Wadler'90],
or $\text{iter}(f)$ in [Trinder'91],
or $\text{ext}(f)$ in [Tannen & Buneman'91].

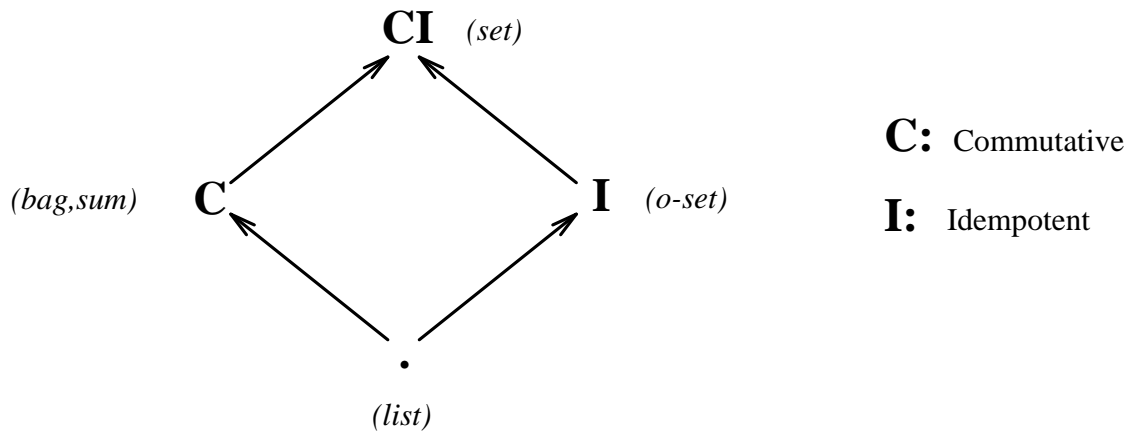
For example,

$$\text{filter}(p)x = \text{hom}[\textit{set}, \textit{set}](\lambda a. \mathbf{if} p(a) \mathbf{then} \{a\} \mathbf{else} \{\}) (x)$$

$$e \in x = \text{hom}[\textit{set}, \textit{some}](\lambda a. (a = e)) (x)$$

$$\text{length}(x) = \text{hom}[\textit{list}, \textit{sum}](\lambda a. 1) (x)$$

Restriction



$\text{hom}[T, S]$ is valid if and only if $T \preceq S$.

$$list \preceq bag \preceq set$$

Valid homomorphism:

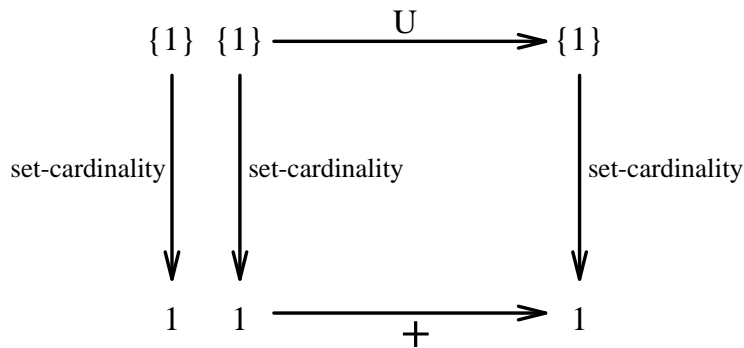
$$\text{bag-cardinality}(x) = \text{hom}[bag, sum](\lambda a.1)(x)$$

since $bag \preceq sum$.

NOT a valid homomorphism:

$$\text{set-cardinality}(x) = \text{hom}[set, sum](\lambda a.1)(x)$$

since $set \not\preceq sum$.



Monoid Comprehensions

For a monoid S :

$$S\{ e \mid r_1, \dots, r_n \} \quad \text{where } r_i \text{ is either } x_i \leftarrow e_i \text{ or } \mathit{pred}$$

$$S\{ e \mid \} \rightarrow \mathit{unit}[S](e)$$

$$S\{ e \mid x \leftarrow u, r_1, \dots, r_n \} \rightarrow \mathit{hom}[T, S](\lambda x. S\{ e \mid r_1, \dots, r_n \})(u)$$

(where $u : T$ and $T \preceq S$)

$$S\{ e \mid \mathit{pred}, r_1, \dots, r_n \} \rightarrow \mathbf{if} \mathit{pred} \mathbf{then} S\{ e \mid r_1, \dots, r_n \} \mathbf{else} \mathit{zero}[S]$$

Example:

$$\begin{aligned} \mathit{join}_{\mathit{pred}}(x, y) &= \mathit{set}\{ (a, b) \mid a \leftarrow x, b \leftarrow y, \mathit{pred} \} \\ &= \mathit{hom}[\mathit{set}, \mathit{set}](\lambda a. \mathit{hom}[\mathit{set}, \mathit{set}](\lambda b. \mathbf{if} \mathit{pred} \mathbf{then} \{ (a, b) \} \\ &\quad \mathbf{else} \{ \}) (y)) (x) \end{aligned}$$

Examples from NF^2 :

Unnesting:

$$\mathit{unnest}(x) = \mathit{set}\{ e \mid s \leftarrow x, e \leftarrow s \}$$

Nesting:

$$\mathit{nest}(k) x = \mathit{set}\{ \langle \mathit{key} : k(e), \mathit{partition} : \mathit{set}\{ a \mid a \leftarrow x, k(e) = k(a) \} \rangle \mid e \leftarrow x \}$$

But you can also join a list with a bag and return a set:

$$\begin{aligned} & \text{set}\{ (a, b) \mid a \leftarrow [1, 2, 3], b \leftarrow \{4, 5\} \} \\ &= \{(1, 4), (1, 5), (2, 4), (2, 5), (3, 4), (3, 5)\} \end{aligned}$$

Other examples:

$$\begin{aligned} \text{filter}(p)(x) &= \text{set}\{ e \mid e \leftarrow x, p(e) \} && x : \text{set} \\ x \cap y &= \text{set}\{ e \mid e \leftarrow x, e \in y \} && x, y : \text{set} \\ \text{length}(x) &= \text{sum}\{ 1 \mid e \leftarrow x \} && x : \text{list} \\ \text{sum}(x) &= \text{sum}\{ e \mid e \leftarrow x \} && x : \text{bag} \\ \text{count}(x, a) &= \text{sum}\{ 1 \mid e \leftarrow x, e = a \} && x : \text{bag} \\ a \in x &= \text{some}\{ a = e \mid e \leftarrow x \} && x : \text{set} \\ \exists a \in x : e &= \text{some}\{ e \mid a \leftarrow x \} && x : \text{set} \\ \forall a \in x : e &= \text{all}\{ e \mid a \leftarrow x \} && x : \text{set} \end{aligned}$$

Translating OQL into the Calculus

select e from x_1 in e_1, \dots, x_n in e_n where $pred$	$bag\{ e \mid x_1 \leftarrow e_1, \dots, x_n \leftarrow e_n, pred \}$
select distinct e from x_1 in e_1, \dots, x_n in e_n where $pred$	$set\{ e \mid x_1 \leftarrow e_1, \dots, x_n \leftarrow e_n, pred \}$
e_1 intersect e_2	$set\{ x \mid x \leftarrow e_1, x \text{ in } e_2 \}$
for all x in e : $pred$	$all\{ pred \mid x \leftarrow e \}$
exists x in e : $pred$	$some\{ pred \mid x \leftarrow e \}$
e_1 in e_2	$some\{ x = e_1 \mid x \leftarrow e_2 \}$
$count(e)$	$sum\{ 1 \mid x \leftarrow e \}$
$sum(e)$	$sum\{ x \mid x \leftarrow e \}$
$flatten(e)$	$T\{ x \mid s \leftarrow e, x \leftarrow s \}$ where $T \in \{set, bag, list\}$
group x in e by $(p_1 : e_1)$	$set\{ \langle p_1 : e_1, partition : set\{ a \mid a \leftarrow e, p_1(a) = p_1(x) \} \rangle \mid x \leftarrow e \}$

Normalization

Canonical forms:

- 1) $\text{hom}[T, S](f) (\textit{path})$
- 2) **if** \textit{path} **then** e_1 **else** e_2
- 3) **if** $(e_1 \otimes e_2)$ **then** e_3 **else** e_4 where $\otimes \in \{=, <, >, \leq, \geq\}$

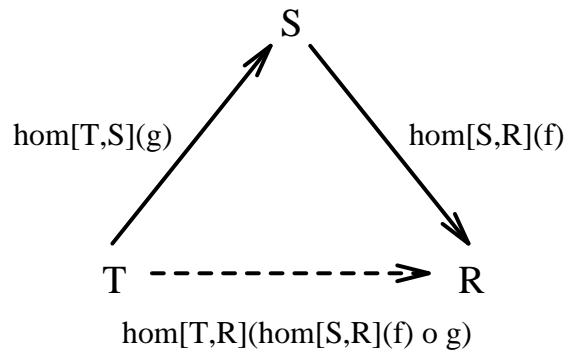
where

$$\begin{aligned} \textit{path} &= \textit{var} \\ \textit{path} &= \textit{path}.\textit{var} \end{aligned}$$

Or equivalently [Wong'93]:

$$T\{ e \mid \dots, x \leftarrow \textit{path}, \dots, e_1 \otimes e_2, \dots \}$$

Normalization rule:



$$\text{hom}[S, R](f) (\text{hom}[T, S](g) (e)) \rightarrow \text{hom}[T, R](\lambda x. \text{hom}[S, R](f) (g(x))) (e)$$

captures the transformation [Wong'93]:

$$T\{ e_1 \mid a \leftarrow S\{ e_2 \mid b \leftarrow x, \bar{r}_1 \}, \bar{r}_2 \} \rightarrow T\{ e_1 \mid b \leftarrow x, \bar{r}_1, a \leftarrow e_2, \bar{r}_2 \}$$

Algebraic optimization:

$\text{filter}(p)(\text{filter}(q)(x))$

$\rightarrow \text{hom}[\text{set}, \text{set}](\lambda a. \mathbf{if } p(a) \mathbf{ then } \{a\} \mathbf{ else } \{\})$
 $\quad (\text{hom}[\text{set}, \text{set}](\lambda a. \mathbf{if } p(a) \mathbf{ then } \{a\} \mathbf{ else } \{\})(x))$

$\rightarrow \text{hom}[\text{set}, \text{set}](\lambda a. \mathbf{if } p(a) \mathbf{ then } (\mathbf{if } q(a) \mathbf{ then } \{a\} \mathbf{ else } \{\}) \mathbf{ else } \{\})(x)$

$= \text{filter}(p \wedge q) x$

Nested OQL query:

select r

from r **in** R

where $r.B$ **in** (**select** $s.D$

from s **in** S

where $r.C = s.C$)

$= \text{set}\{r \mid r \leftarrow R, \text{some}\{x = r.B \mid x \leftarrow \text{set}\{s.D \mid s \leftarrow S, r.C = s.C\}\}\}$

$\rightarrow \text{set}\{r \mid r \leftarrow R, s \leftarrow S, s.D = r.B, r.C = s.C\}$

Why canonical forms are important?

- has been proved so for other domains (e.g. disjunctive normal forms);
- it is easier to prove general theorems about properties of programs in canonical form [Wong'93];
- normalization makes programs more efficient in *most* cases;
- the evaluation of canonical forms generates less intermediate data structures;
- normalized programs are amenable to a higher degree of parallelism.

Current extensions

The algebra can be easily extended to support:

- arrays with *real* matrix manipulations;
- object identity (yet another monoid);
- destructive database updates;
- non-deterministic operations (e.g. *listify*);
- pattern matching in comprehensions.