

Comprehending Monadic Queries

Jeremy Gibbons (joint work with Fritz Henglein, Ralf Hinze, Nicolas Wu) WG2.11#15, November 2015

1. Comprehensions

• ZF axiom schema of specification:

$$\{x^2 \mid x \in Nat \land x < 10 \land x \text{ is even}\}$$

• SETL set-formers:

```
\{x * x : x \text{ in } \{0...9\} \mid x \text{ mod } 2 = 1\}
```

• Eindhoven Quantifier Notation:

$$(x:0 \le x < 10 \land x \text{ is even}: x^2)$$

• Haskell (NPL, Python, ...) list comprehensions:

$$[x \land 2 \mid x \leftarrow [0..9], even x]$$

2. Relational algebra vs calculus

Consider two database tables:

```
customers: <u>cid</u>, name, address
invoices : iid, customer, amount, due
```

A query in relational *algebra* ('point-free', on relations):

```
\pi_{name,amount,address} (\sigma_{due < today} (customers \bowtie_{cid = customer} invoices))
```

The same query in relational *calculus* ('point-wise', on tuples):

```
SELECT name, amount, address
FROM customers, invoices
WHERE cid = customer AND due < today
```

The algebraic style may be convenient for formal manipulation, but the calculus style is much more accessible for readers. DBMSs typically translate from calculus-style input to algebra-style intermediate representation.

3. Comprehending queries

Trinder (1991) argued for comprehensions as a query notation:

```
[ (c.name, c.address, i.amount)
| c ← customers,
i ← invoices,
c.cid == i.customer,
i.due < today]</pre>
```

Very influential observation in the DBPL community. Formed the basis of languages such as Buneman's *Kleisli*, Microsoft *LINQ*, Wadler's *Links*, as well as querying for objects (*OQL*) and XML (*XQuery*).

4. Comprehending monads (Wadler 1992)

The necessary structure is that of a *monad* $(T, \gg, return)$:

(>=) ::
$$T \ a \rightarrow (a \rightarrow T \ b) \rightarrow T \ b$$
 $(x \gg f) \gg k = x \gg (\lambda a \rightarrow f \ a \gg k)$
return :: $a \rightarrow T \ a$ return $a \gg k = k \ a$
 $x \gg return = x$

with additionally *mzero* :: T *a*.

Comprehensions can then be generalized to other monads:

```
\mathcal{D} [e|] = return e
\mathcal{D} [e|p \leftarrow e', Q] = e' \gg \lambda p \rightarrow \mathcal{D} [e|Q]
\mathcal{D} [e|e', Q] = guard e' \gg \mathcal{D} [e|Q]
\mathcal{D} [e|\mathbf{let} d, Q] = \mathbf{let} d \mathbf{in} \mathcal{D} [e|Q]
```

(where guard b = if b then return () else mzero).

Hence monad comprehensions for sets, bags, maps-to-monad-zeroes, etc.

5. The problem with joins

The comprehension yields a terrible query plan!

Constructs entire cartesian product, then discards most of it:

```
cp customers invoices \triangleright filter (\lambda(c,i) \rightarrow c.cid == i.customer) \triangleright filter (\lambda(c,i) \rightarrow i.due < today) \triangleright fmap (\lambda(c,i) \rightarrow (c.name, c.address, i.amount)
```

(where \triangleright is reverse function application).

Better to group by customer identifier, then handle groups separately:

```
(indexBy cid customers) 'merge' (indexBy customer invoices) \triangleright fmap (id \times filter (\lambda i \rightarrow i.due < today)) \triangleright fmap (fmap (\lambda c \rightarrow (c.name, c.address)) <math>\times fmap (\lambda i \rightarrow i.amount))
```

(where *indexBy* partitions, and *merge* pairs on common index). But this doesn't correspond to anything expressible in comprehensions.

6. Comprehensive comprehensions

Various extensions to the comprehension syntax:

• parallel ('zip') comprehensions (since GHC 5.0, 2001):

```
[(x,y) | x \leftarrow [1,2,3] | y \leftarrow [4,5,6]]
```

• 'order by' and 'group by' (Wadler & Peyton Jones, 2007):

```
[ (the dept, sum salary)
| (name, dept, salary) ← employees
, then group by dept using groupWith
, then sortWith by sum salary
, then take 5]
```

(NB **group by** rebinds the variables bound earlier!)

Initially just for lists, but...

Generalized comprehensive comprehensions

... generalizes nicely to other monads (Giorgidze et al, 2011):

```
\mathcal{D} [e \mid (Q \mid R), S]
= mzip (\mathcal{D} [vQ \mid Q]) (\mathcal{D} [vR \mid R]) \gg \lambda(vQ, vR) \rightarrow \mathcal{D} [e \mid S]
\mathcal{D} [e \mid Q, \mathbf{then} f \mathbf{by} b, R]
= f (\lambda vQ \rightarrow b) (\mathcal{D} [vQ \mid Q]) \gg \lambda vQ \rightarrow \mathcal{D} [e \mid R]
\mathcal{D} [e \mid Q, \mathbf{then} \mathbf{group} \mathbf{by} b \mathbf{using} f, R]
= f (\lambda vQ \rightarrow b) (\mathcal{D} [vQ \mid Q]) \gg \lambda ys \rightarrow
\mathbf{case} (fmap vQ_1 ys, ..., fmap vQ_n ys) \mathbf{of} vQ \rightarrow \mathcal{D} [e \mid R]
```

where vQ is the tuple of variables bound by Q (and used subsequently), and vQ_i is a selector mapping vQ to its ith component.

7. Solving the problem with (equi-)joins

Maps-to-bags form a monad-with-zero—roughly:

```
type Map \ k \ v = k \rightarrow v
type Table \ k \ v = Map \ k \ (Bag \ v)
```

Now define

```
indexBy :: Eq k \Rightarrow (v \rightarrow k) \rightarrow Bag \ v \rightarrow Table \ k \ v
indexBy f xs k = filter \ (\lambda v \rightarrow f \ v == k) xs
merge :: Table k \ v \rightarrow Table \ k \ w \rightarrow Table \ k \ (v, w)
merge f g = \lambda k \rightarrow cp \ (f \ k) \ (g \ k)
```

Can use *merge* for parallel comprehensions:

instance MonadZip (Table k) where mzip = merge and indexBy for grouping.

Given input tables

```
customers:: Bag (CID, Name, Address) invoices:: Bag (IID, CID, Amount, Date)
```

evaluate our example query as:

Avoids expanding the whole cartesian product.

8. Aggregation

For database queries, want to aggregate collections: *count*, *sum*, *some*, . . .

Problem: maps may be infinite.

Solution: restrict to *finite* maps.

Problem: not a monad—*return* $a = \lambda k \rightarrow a$ yields a non-finite map.

Solution? *semi-monads* (with bind but no return).

Problem: semi-monad comprehensions—base case uses *return*:

 $\mathcal{D}[e] = return e$

This is surmountable...but we prefer:

Solution: *graded* (indexed, parametric) monads

9. Graded monads

Monad $(T, \gg, return)$ has endofunctor $T: C \to C$, polymorphic functions

$$(\gg)$$
 :: T $a \rightarrow (a \rightarrow T b) \rightarrow T b$
return :: $a \rightarrow T a$

such that

$$(x \gg f) \gg k = x \gg (\lambda a \rightarrow f \ a \gg k)$$

return $a \gg k = k \ a$
 $x \gg return = x$

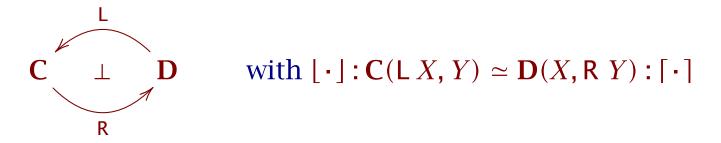
Katsumata's M-graded monad $(\mathsf{T},\gg, return)$ for monoid (M,\cdot,ε) has (non-endo-)functor $\mathsf{T}:M\to [\mathsf{C},\mathsf{C}]$ and

```
(≫=) :: T m a \rightarrow (a \rightarrow T n b) \rightarrow T (m \cdot n) b
return :: a \rightarrow T \varepsilon a
```

with same laws. We use T = Table over monoid $(K, \times, 1)$ of *finite* key types.

10. Adjunctions, and query optimization

Optimizations depend on a body of *meaning-preserving transformations*, all arising from algebraic properties of the datatypes—*adjunctions*:



Currying yields indexing; *products* yield projection and merge; *coproducts* yield filters; *free* commutative monoids yield selection and aggregation.

Monads famously arise from adjunctions; graded monads do too, albeit in a slightly more complicated way.

Work in progress: justifying standard query optimizations via these correspondences.

11. Comprehending semi-monads

Prohibit comprehensions with no qualifiers; multiple base cases instead.

```
\mathcal{D}\left[\varepsilon\mid p\leftarrow e'\right] = fmap\left(\lambda p\rightarrow e'\right)\varepsilon
\mathcal{D}\left[\varepsilon\mid e'\right] = \dots - \text{not allowed}
\mathcal{D}\left[\varepsilon\mid \mathbf{let}\ d\right] = \dots - \text{not allowed}
\mathcal{D}\left[\varepsilon\mid (Q\mid R)\right] = fmap\left(\lambda(vQ, vR) \rightarrow \varepsilon\right)
\left(mzip\left(\mathcal{D}\left[vQ\mid Q\right]\right)\left(\mathcal{D}\left[vR\mid R\right]\right)\right)
\mathcal{D}\left[\varepsilon\mid Q, \mathbf{then}\ f\ \mathbf{by}\ b\right] = fmap\left(\lambda vQ \rightarrow \varepsilon\right)\left(f\left(\lambda vQ \rightarrow b\right)\left(\mathcal{D}\left[vQ\mid Q\right]\right)\right)
\mathcal{D}\left[\varepsilon\mid Q, \mathbf{then}\ \mathbf{group}\ \mathbf{by}\ b\ \mathbf{using}\ f\right]
= fmap\left(\lambda ys \rightarrow \mathbf{case}\ (fmap\ vQ_1\ ys, ..., fmap\ vQ_n\ ys)\ \mathbf{of}\ vQ \rightarrow \varepsilon\right)
\left(f\left(\lambda vQ \rightarrow b\right)\left(\mathcal{D}\left[vQ\mid Q\right]\right)\right)
```

Also, we can't define *guard* if we don't have *return*, so desugaring of guards needs to change:

```
\mathcal{D} [\varepsilon \mid e', Q] = \text{if } e' \text{ then } \mathcal{D} [\varepsilon \mid Q] \text{ else } mzero
```