

An Effective Framework for Processing Object-Oriented Database Languages

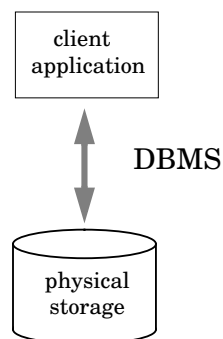
Leonidas Fegaras

Oregon Graduate Institute

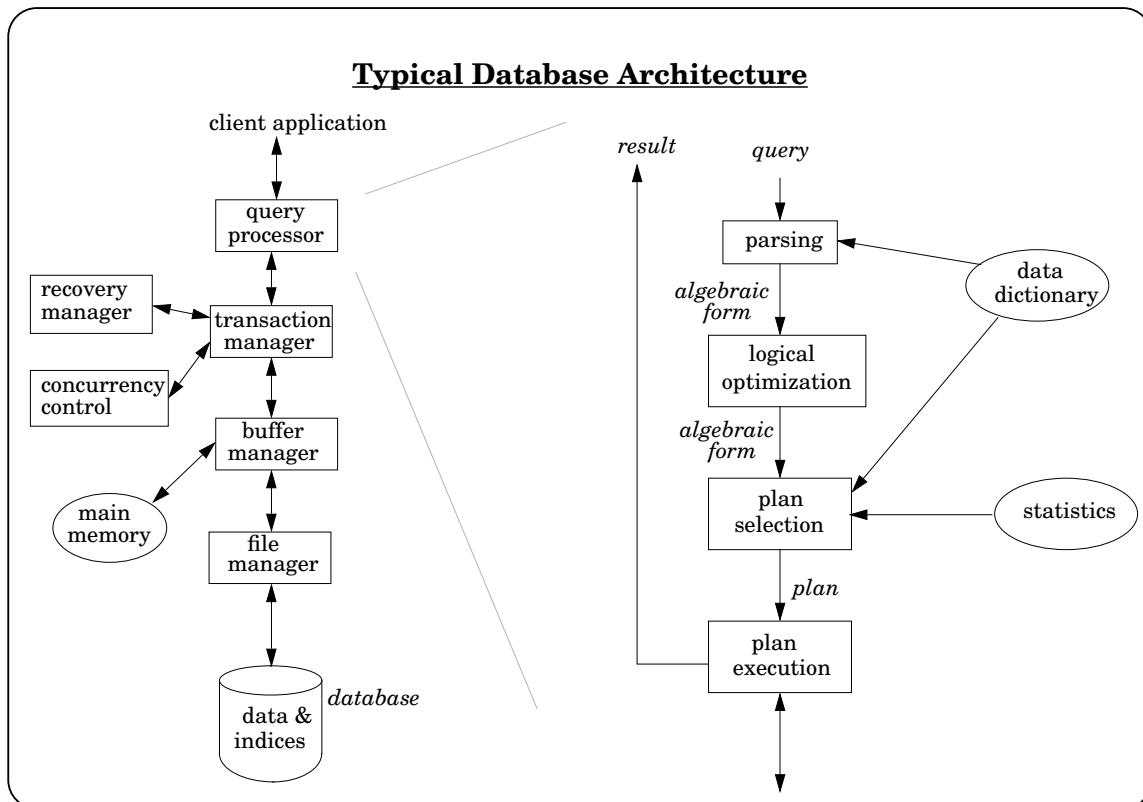
(Ph.D. University of Massachusetts)

Maintaining Persistent Data

Data Base Management System:



Key issue: *data independence*.



The Gap Between Theory & Practice

Most commercial relational query languages are based on

- tuple calculus (SQL, Quel)
- domain calculus (QBE)

However, in some respects they go beyond the formal model:

- aggregate operators,
- sort orders,
- grouping,
- update capabilities.

Insufficient Modeling Power

Relational databases cannot effectively model many new applications:

- multimedia & World-Wide-Web,
- scientific databases,
- CAD,
- CASE,
- GIS,
- office automation.

New Requirements

New database languages must be able to handle:

- type extensibility;
- multiple collection types (e. g., sets, lists, trees, arrays);
- arbitrary nesting of type constructors;
- large objects (e.g., text, sound, image);
- temporal & spatial data;
- unstructured data;
- active rules;
- methods.

New Proposals for Database Languages

Relational extensions:

- UniSQL,
- POSTGRES/Illustra,
- SQL3.

Object-oriented databases:

- Gemstone,
- O2,
- OQL of ODMG-93.

Why Do We Need a Formal Calculus?

- facilitates equational reasoning;
- provides a theory for proving query transformations correct;
- imposes language uniformity;
- avoids language inconsistencies.

functional languages	\longleftrightarrow	lambda calculus
relational databases	\longleftrightarrow	relational algebra/calculus
object-oriented databases	\longleftrightarrow	?

What is an Effective Calculus?

Several aspects:

- coverage,
- ease of manipulation,
- ease of evaluation,
- uniformity.

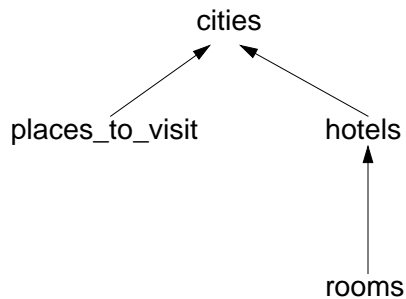
Rest of the Talk

- Monoids;
- *Algebra*: monoid homomorphisms;
- *Calculus*: monoid comprehensions;
- unnesting comprehensions;
- overview of the query translation framework;
- physical design specification;
- one example of query translation.

Case Study: ODMG-93 OQL

```
class city = < name: string,  
             hotels: bag (hotel),  
             places_to_visit: list (< name: string, address: string >) >  
extent cities;
```

```
class hotel = < name: string,  
              rooms: set (< bed#: int, price: int >) >;
```



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

OQL:

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

Monoid comprehension:

```
set { h.name | c ← cities,  
             h ← c.hotels,  
             p ← c.places_to_visit,  
             c.name="Portland",  
             h.name=p.name }
```

Monoids

A *monoid* is an algebraic structure that captures most collection and aggregate types:

<i>operator</i>	<i>functionality</i>	e.g., sets
zero	the identity value	{ }
merge(x,y)	associative with identity zero	$x \cup y$
unit(a)	singleton construction	{ a }

$$\{1, 2, 3\} = \{1\} \cup \{2\} \cup \{3\}$$

Optional properties:

commutativity : **merge(x,y) = merge(y,x)**

idempotence : **merge(x,x) = x**

Some Monoids

Collection Monoids

monoid	type	zero	unit(a)	merge
list	list(α)	[]	[a]	append
set	set(α)	{ }	{a}	\cup
bag	bag(α)	{ { }	{ {a} }	$\dot{\cup}$
sorted[f]	list(α)	[]	[a]	list_merge[f]

Primitive Monoids

monoid	type	zero	unit(a)	merge
sum	integer	0	a	+
prod	integer	1	a	*
some	boolean	false	a	\vee
all	boolean	true	a	\wedge

Monoid Homomorphisms

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

$$\alpha \xrightarrow{f} S \quad T(\alpha) \xrightarrow{\mathbf{hom}[T,S](f)} S$$

$$\mathbf{hom}[T,S](f)(\mathbf{zero}[T]) = \mathbf{zero}[S]$$

$$\mathbf{hom}[T,S](f)(\mathbf{unit}[T](a)) = f(a)$$

$$\mathbf{hom}[T,S](f)(\mathbf{merge}[T](x,y)) = \mathbf{merge}[S](\mathbf{hom}[T,S](f)(x), \mathbf{hom}[T,S](f)(y))$$

Example

Notation: $\lambda a . \{ a^2 \}$ is the function f such that: $f(a) = \{ a^2 \}$

$$\begin{aligned} & \mathbf{hom}[bag, set](\lambda a . \{ a^2 \}) \{ \{ 1 \}, \{ 2 \}, \{ 3 \} \} \\ &= \mathbf{hom}[bag, set](f) (\{ \{ 1 \} \} \uplus \{ \{ 2 \} \} \uplus \{ \{ 3 \} \}) \\ &= f(1) \cup f(2) \cup f(3) \\ &= \{ 1 \} \cup \{ 4 \} \cup \{ 9 \} \\ &= \{ 1, 4, 9 \} \end{aligned}$$

Other Examples

$$5 \in x = \mathbf{hom}[\mathit{set}, \mathit{some}](\lambda a. (a=5))(x)$$

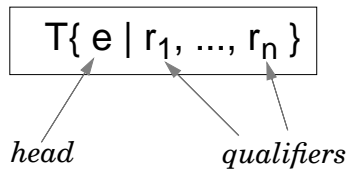
$$\begin{array}{c} \{8\} \cup \{5\} = \{8, 5\} \\ \downarrow \quad \downarrow \\ 8=5 \quad 5=5 \\ \downarrow \quad \downarrow \\ \text{false} \vee \text{true} = \text{true} \end{array}$$

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

$$\begin{array}{c} \mathit{append}([2], [3]) = [2, 3] \\ \downarrow \quad \downarrow \\ 1 \quad 1 \\ \downarrow \quad \downarrow \\ 1 + 1 = 2 \end{array}$$

Monoid Comprehensions

A *monoid comprehension* takes the form:



where T is a monoid and each *qualifier* r_i is either:

- a *generator* $v \leftarrow u$;
- a *filter* pred.

Examples

$$\text{set} \{ (a,b) \mid a \leftarrow x, b \leftarrow y \} = \left\{ \begin{array}{l} \text{res} = \{ \}; \\ \text{for each } a \text{ in } x \text{ do} \\ \quad \text{for each } b \text{ in } y \text{ do} \\ \quad \quad \text{res} = \text{res} \cup \{ (a,b) \}; \\ \text{return res;} \end{array} \right.$$

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

$$\text{sum} \{ a \mid a \leftarrow [1,2,3], a \geq 2 \} = 2+3 = 5$$

Formal Definition of a Monoid Comprehension

$$T\{ e \mid \} = \text{unit}[T](e)$$

$$T\{ e \mid v \leftarrow u, r_1, \dots, r_n \} = \text{hom}[S,T](\lambda v. T\{ e \mid r_1, \dots, r_n \})(u)$$

where S is the type of u

$$T\{ e \mid \text{pred}, r_1, \dots, r_n \} = \text{if pred then } T\{ e \mid r_1, \dots, r_n \} \text{ else zero}[T]$$

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

$$= \text{hom}[\text{list}, \text{set}](\lambda a. \text{hom}[\text{bag}, \text{set}](\lambda b. \{ (a,b) \})$$

$$\quad \quad \quad (\{ \{4,5\} \})$$

$$\quad \quad \quad ([1,2,3])$$

Other Examples

$\text{filter}(\text{pred}) e = \text{set}\{ x \mid x \leftarrow e, \text{pred}(x) \}$
 $\text{flatten}(e) = \text{set}\{ x \mid s \leftarrow e, x \leftarrow s \}$
 $e_1 \cap e_2 = \text{set}\{ x \mid x \leftarrow e_1, x \in e_2 \}$
 $\text{length}(e) = \text{sum}\{ 1 \mid x \leftarrow e \}$
 $\exists a \in e: \text{pred} = \text{some}\{ \text{pred} \mid a \leftarrow e \}$
 $\forall a \in e: \text{pred} = \text{all}\{ \text{pred} \mid a \leftarrow e \}$
 $\text{nest}(k) e = \text{set}\{ \langle \text{KEY} = k(x), \text{DATA} = \text{set}\{ y \mid y \leftarrow e, k(x) = k(y) \} \rangle \mid x \leftarrow e \}$
 $\text{unnest}(e) = \text{set}\{ x \mid s \leftarrow e, x \leftarrow s.\text{DATA} \}$

Translating OQL

```
select distinct h.name  
from hl in ( select c.hotels  
             from c in cities  
             where c.name="Portland" ),  
  h in hl  
where exists r in h.rooms: ( r.bed#=3 )
```

```
set{ h.name | hl ← bag{ c.hotels | c ← cities, c.name="Portland" },  
      h ← hl,  
      some{ r.bed#=3 | r ← h.rooms } }
```

Program Normalization

Canonical form: (path is a cascade of projections: $X.A_1.A_2 \dots A_m$)

- $T\{ e \mid x_1 \leftarrow \text{path}_1, \dots, x_n \leftarrow \text{path}_n, \text{pred} \}$

Examples of normalization rules:

- $T\{ e \mid \boxed{\textcircled{1}}, x \leftarrow S\{ u \mid \boxed{\textcircled{2}} \}, \boxed{\textcircled{3}} \}$
 $\rightarrow T\{ e \mid \boxed{\textcircled{1}}, \boxed{\textcircled{2}}, x \equiv u, \boxed{\textcircled{3}} \}$
- $T\{ e \mid \boxed{\textcircled{1}}, \text{some}\{ \text{pred} \mid \boxed{\textcircled{2}} \}, \boxed{\textcircled{3}} \}$
 $\rightarrow T\{ e \mid \boxed{\textcircled{1}}, \boxed{\textcircled{2}}, \text{pred}, \boxed{\textcircled{3}} \}$

Example

$\text{set}\{ h.\text{name} \mid h1 \leftarrow \text{bag}\{ c.\text{hotels} \mid c \leftarrow \text{cities}, c.\text{name} = \text{"Portland"} \},$
 $h \leftarrow h1,$
 $\text{some}\{ r.\text{bed}\# = 3 \mid r \leftarrow h.\text{rooms} \} \}$

= $\text{set}\{ h.\text{name} \mid c \leftarrow \text{cities}, c.\text{name} = \text{"Portland"},$
 $h1 \equiv c.\text{hotels},$
 $h \leftarrow \underline{h1},$ Substitute c.hotels for h1
 $\text{some}\{ r.\text{bed}\# = 3 \mid r \leftarrow h.\text{rooms} \} \}$

= $\text{set}\{ h.\text{name} \mid c \leftarrow \text{cities}, c.\text{name} = \text{"Portland"},$
 $h \leftarrow c.\text{hotels},$
 $\text{some}\{ r.\text{bed}\# = 3 \mid r \leftarrow h.\text{rooms} \} \}$

= $\text{set}\{ h.\text{name} \mid c \leftarrow \text{cities},$
 $h \leftarrow c.\text{hotels},$
 $r \leftarrow h.\text{rooms},$
 $(c.\text{name} = \text{"Portland"}) \wedge (r.\text{bed}\# = 3) \}$

Unnesting OQL Queries

```
select distinct h.name
from hl in ( select c.hotels
             from c in cities
             where c.name="Portland" ),
h in hl
where exists r in h.rooms: ( r.bed#=3 )
```



```
select distinct h.name
from c in cities,
h in c.hotels,
r in h.rooms
where c.name="Portland" and r.bed#=3
```

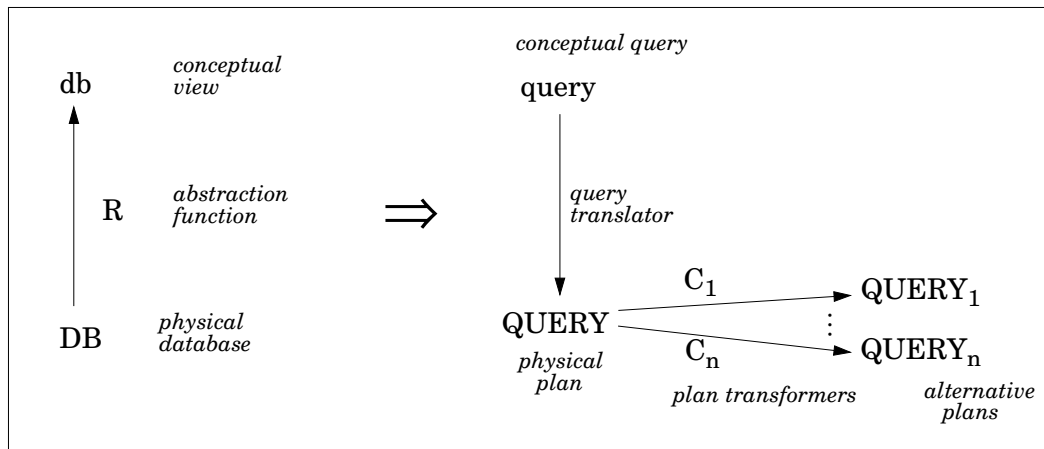
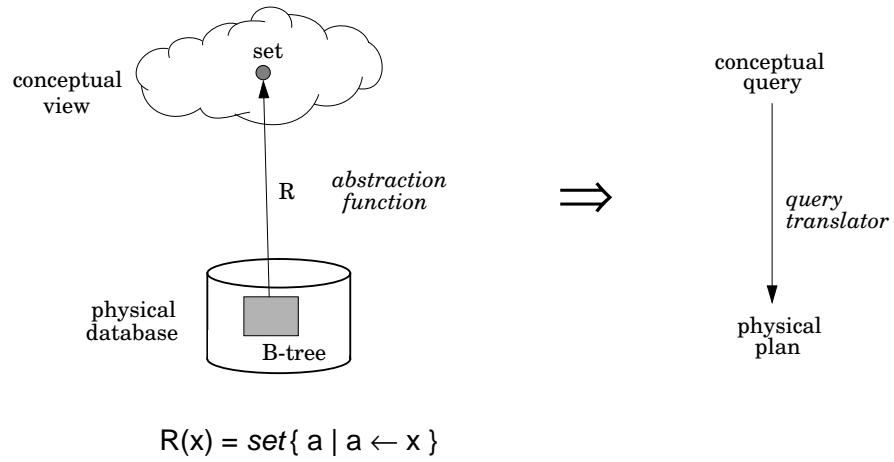
Query Translation

Need to support data independence.

Tasks:

- the *database administrator* specifies the conceptual database schema;
- the *database implementor* specifies the *physical design*;
- an *application programmer* submits a query against the database;
- the *query translator* translates the query into a physical plan that reflects the physical design.

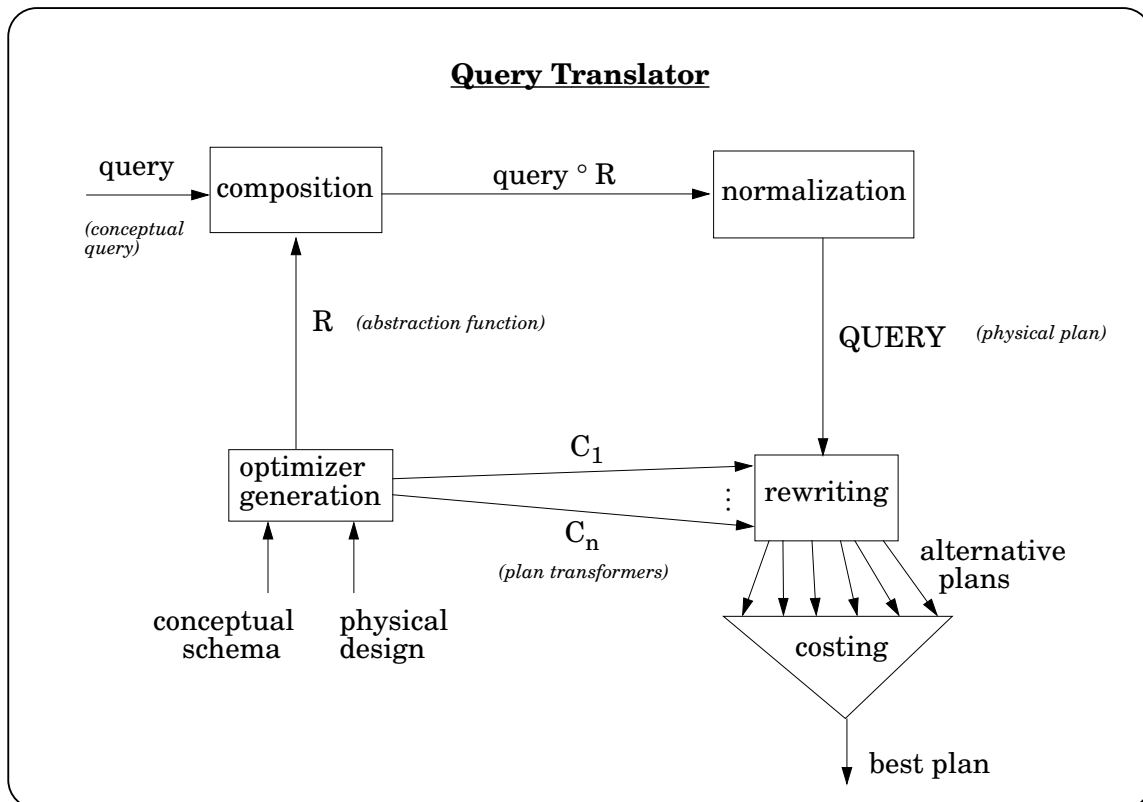
Conceptual-to-Internal Mapping



$$\begin{aligned} \text{db} &= R(\text{DB}) \\ \text{QUERY}(\text{DB}) &= \text{query}(R(\text{DB})) && \text{physical plan} \end{aligned}$$

Plan transformers:

$$\begin{aligned} \text{DB} &= C_i(\text{DB}) \\ \text{QUERY}_i(\text{DB}) &= \text{QUERY}(C_i(\text{DB})) && \text{alternative plans} \end{aligned}$$

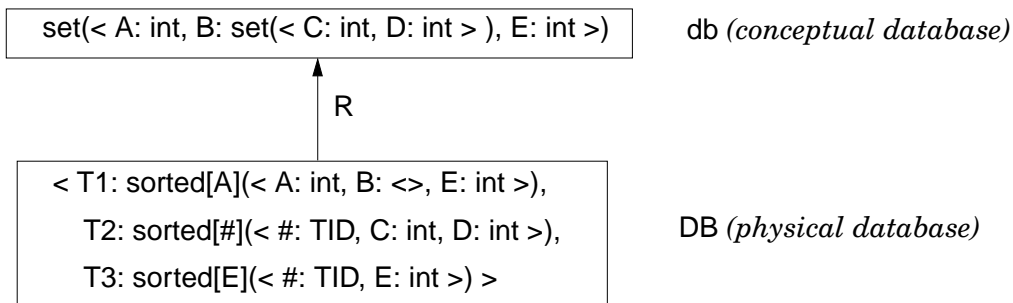


Example

Conceptual Schema:
 $\text{set}(\langle A: \text{int}, B: \text{set}(\langle C: \text{int}, D: \text{int} \rangle), E: \text{int} \rangle)$

Physical design:

- normalize the inner set;
- implement the outer set as a B-tree with key A;
- attach a secondary index to the outer set with a key E.



Abstraction function:

$$R(DB) = \text{set}\{ \langle A = a.A, \\ B = \text{set}\{ \langle C = b.C, D = b.D \rangle \mid b \leftarrow DB.T2, b.\# = @a \}, \\ E = a.E \rangle \\ \mid a \leftarrow DB.T1 \}$$

Plan Transformer: $DB = C_1(DB)$ *(reconstructs T1 from T3)*

Conceptual query:

$$\text{query}(db) = \text{sum}\{ y.C \mid x \leftarrow db, y \leftarrow x.B, x.A=10, y.D>5 \}$$

Physical plan:

$$\begin{aligned} \text{QUERY}(DB) &= \text{query}(R(DB)) \\ &= \text{sum}\{ y.C \mid x \leftarrow R(DB), y \leftarrow x.B, x.A=10, y.D>5 \} \\ &= \dots \quad (\text{after normalization}) \\ &= \text{sum}\{ b.C \mid a \leftarrow DB.T1, b \leftarrow DB.T2, \\ &\quad b.\# = @a, a.A=10, b.D>5 \} \end{aligned}$$

(A sort-merge join!)

Alternative plan:

$$\begin{aligned} \text{QUERY}_1(DB) &= \text{QUERY}(C_1(DB)) \\ &= \dots \quad (\text{a plan that uses the secondary index T3}) \end{aligned}$$

Our Query Translation Framework is Purely Algebraic

It formalizes and generalizes many ad-hoc techniques used in relational query translation and optimization:

- the mapping from logical queries to physical plans is *meaning preserving*;
- the data translation overhead is completely eliminated by normalization;
- alternative physical plans are derived by a very simple cost-directed term-rewriting system that uses 3 types of rules:
 - associativity of monoids;
 - commutativity of some monoids;
 - plan transformers.

Physical Design Specification

A *physical design language* is provided that is

- declarative,
- extensible.

Captures many physical designs:

- object clustering;
- horizontal/vertical partitioning;
- schema normalization;
- join indices;
- multiple access paths via secondary indices.

Conclusion

I have presented:

- a uniform calculus that captures many new database language features:
 - supports expression nesting;
 - allows arbitrary nesting of type constructors;
 - supports multiple collection types;
 - handles aggregates & predicates directly.
- a normalization algorithm that unnests nested comprehensions;
- an effective algebraic model for query translation that facilitates data independence.

Current and Future Work

Model extensions:

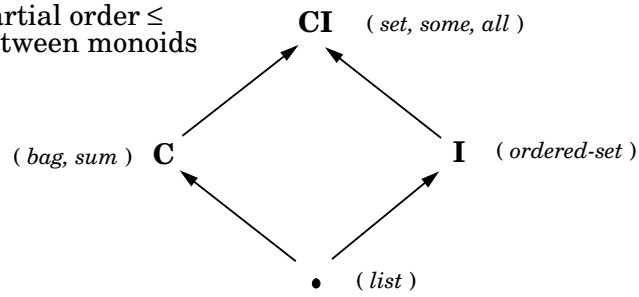
- vectors and arrays;
- object identity & mutable objects;
- updates;
- methods.

Implementation:

- translators from OQL and a subset of SQL3 to the monoid calculus;
- a query optimizer.

Restriction

Partial order \leq
between monoids

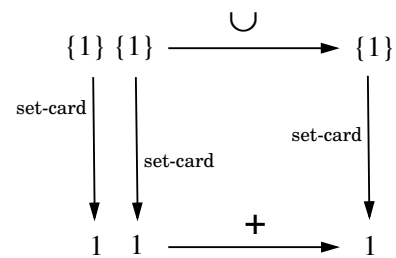


C: Commutative
I: Idempotent

$\mathbf{hom}[T,S]$ is *well-formed* if and only if $T \leq S$

bag-cardinality(x) = $\mathbf{hom}[bag, sum](\lambda a.1)(x)$
is well-formed, since $bag \leq sum$

set-cardinality(x) = $\mathbf{hom}[set, sum](\lambda a.1)(x)$
is **not** well-formed, since $set > sum$

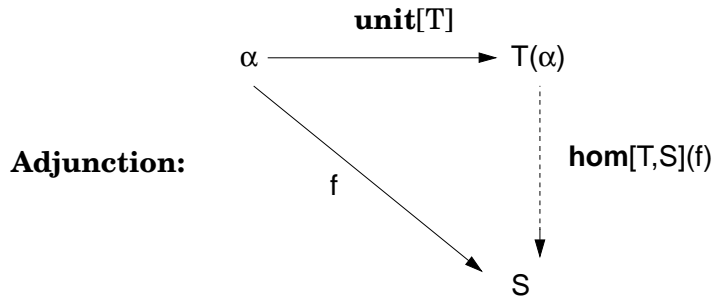


Related Work

- *monoid homomorphisms* [V. Tannen et al];
 - SRU
 - $\text{ext}(f) A = \mathbf{hom}[T,T](f) A$
- *boom hierarchy of types* [R. Bird, L. Meertens, R. Backhouse];
- *monad comprehensions* [P. Wadler, P. Trinder, P. Buneman];
- *normalization* [L. Wong].

In Terms of Category Theory

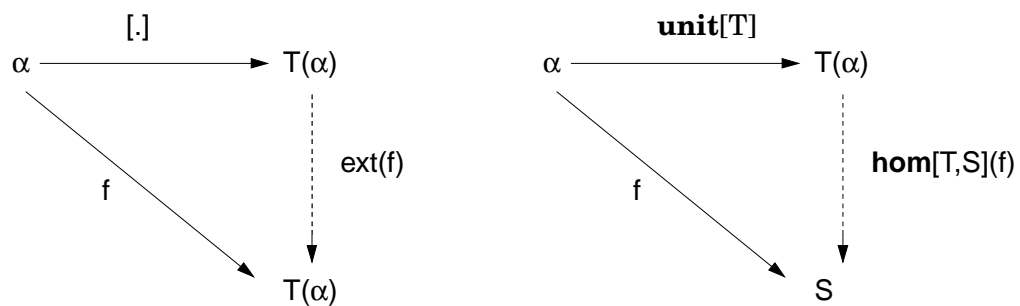
The collection monoid $T(\alpha)$ is a **free monoid** generated by α :



$\mathbf{unit}[T]$ is a natural transformation: $\text{id} \rightarrow T$.

$\mathbf{hom}[T,S](f)$ is the **homomorphic extension** of f .

Monads vs. Monoids



The Kleisli triple $(T, [.] , \text{ext}(f))$, which represents a monad, is a monoid:

the Kleisli composition $f \bullet g = \text{ext}(f) \circ g$ is associative with zero $[.]$.