# Investigation of Algebraic Query Optimisation for Database Programming Languages

Alexandra Poulovassilis
Dept. of Computer Science
King's College London
Strand, London WC2R 2LS
alex@dcs.kcl.ac.uk

Carol Small
Dept. of Computer Science
Birkbeck College
Malet St., London WC1E 7HX
carol@dcs.bbk.ac.uk

## Abstract

A major challenge still facing the designers
and implementors of database programming
languages (DBPLs) is that of query optimisa-
tion. We investigate algebraic query optimi-
sation techniques for DBPLs in the context of
a purely declarative functional language that
supports sets as first-class objects. Since the
language is computationally complete issues
such as non-termination of expressions and
construction of infinite data structures can be
investigated, whilst its declarative nature al-
lows the issue of side effects to be avoided and
a richer set of equivalences to be developed.
The support of a set bulk data type enables
much prior work on the optimisation of rela-
tional languages to be utilised. Finally, the
language has a well-defined semantics which
permits us to reason formally about the prop-
erties of expressions, such as their equivalence
with other expressions and their termination.

## 1 Introduction

Database programming languages (DBPLs) incorpo-
rate into a single language, with a single semantics, all
of the features normally expected of both a data ma-
nipulation language (DML) and a programming lan-
guage. For example, DBPLs have one computational
model, one type system, and bulk data types with as-
sociated access mechanisms. A major challenge still
facing DBPLs is that of query optimisation. There
are several reasons for limited progress in this area:

(i) The possibility of *side-effects* restricts the set of
equivalences that can be shown to hold.

(ii) Some bulk data structures are inherently hard to
optimise. For example, lists only readily support
the common relational optimisations if the con-
cept of *bag equality* is used [Tri89] (i.e. lists are
equal if they contain the same elements, although
possibly in different orders), whilst some of the
algebraic properties of sets fail for bags [Alb91].

(iii) Since DBPLs are computationally complete, the
*termination* properties of expressions must be
taken into account when investigating equiva-
lences. For example, if the boolean-valued func-
tion $f_1$ does not terminate for some arguments
whilst the boolean-valued function $f_2$ returns
*False* for all arguments, then the 'equivalence'
$\sigma_{f_1}(\sigma_{f_2}(s)) = \sigma_{f_2}(\sigma_{f_1}(s))$ does not hold since
evaluation of the LHS always terminates for finite
$s$ (returning $\{\}$) whereas evaluation of the RHS
may not terminate.

(iv) DBPLs may manipulate *infinite data structures*
and hence some bulk data operations cannot be
implemented using established methods. For ex-
ample, if $A$ and $B$ are infinite sets then a nested
loop method cannot be used to generate $A \times B$
(since all tuples of the resulting product will have
the same first coordinate).

(v) DBPLs typically support *user-defined data types*,
and hence require mechanisms to prove equiva-

lences over these data types too.

In this paper we investigate optimisation techniques for DBPLs by addressing some of the above issues. We undertake our investigation in the context of a purely declarative functional language. Since database algebras are functional in nature, this is a particularly natural computational paradigm in which to investigate query optimisation in DBPLs. It also gives us a computationally complete formalism that can exhibit non-termination of expressions (point (iii) above) and that can result in infinite data structures (point (iv) above), whilst avoiding the issue of side-effects (point (i) above). The language supports a set bulk data type, enabling us to utilise much prior work on the optimisation of relational languages, including Datalog (point (ii) above). Finally, the language has a well-defined semantics which permits us to reason formally about the properties of expressions, including those of user-defined data types, such as their equivalence with other expressions and their termination properties (point (v) above).

The structure of this paper is as follows. In Section 2 we give the syntax of our language and briefly discuss its semantics and its provision for built-in and user-defined functions. In Section 3 we define a small but powerful algebra of operators over the set data type, provide some key equivalences for expressions in these operators, and list transformation principles for optimising expressions. Along the way, we identify some caveats to well-known equivalences for non-deductive database languages. In Section 4 we examine two higher level constructs commonly found in DBPLs - *set abstractions* (also known as *set comprehensions* in the literature), and functions with inverses. We provide some key equivalences for these constructs also, and give transformation principles for expressions in them. In Section 5 we briefly compare this work with related research. Finally, in Section 6 we give our conclusions and indicate directions of further work.

## 2 The Language

The formal foundation of any functional language is the $\lambda$-calculus [Hin86]. Expressions in this calculus have the following syntax:

$$expr \quad = \quad var \mid primitive \mid \text{``}\lambda\text{''}var\text{``.''}expr \mid$$
$$expr_1 \ expr_2 \mid \text{``(''} \ expr \ \text{``)''}$$

A variable $x$ is said to be *bound* in an expression $e$ if it occurs in a sub-expression of $e$ of the form $\lambda x.e'$; otherwise it is *free* in $e$. $FV(e)$ (respectively, $BV(e)$) denotes the set of variables with at least one free (respectively, bound) occurrence in $e$.

Computation in the $\lambda$ calculus proceeds by syntactically transforming terms using $\beta$ reduction. This rewrites a function application of the form $(\lambda x.e)e'$ to the expression $e[e'/x]$ obtained by replacing all free occurrences of $x$ in $e$ by $e'$. The denotational semantics of the $\lambda$-calculus (see [Sch86]) assigns to each expression a value in a semantic domain - this is the meaning of the expression. $\beta$ reduction is semantically sound in that it does not alter the meaning of an expression.

The language that we will be optimising is the $\lambda$-calculus extended with constructors, *let* expressions and pattern-matching $\lambda$-abstractions:

$$expr \quad = \quad var \mid constructor \mid primitive \mid$$
$$\text{``}\lambda\text{''}pattern\text{``.''}expr \mid expr_1 \ expr_2 \mid$$
$$\text{``let''} \ var \ \text{``=''} \ expr_1 \ \text{``in''} \ expr_2 \mid$$
$$\text{``(''} \ expr \ \text{``)''}$$
$$pattern = \quad var \mid constructor \ pattern_1 \ ... \ pattern_n$$

where tuples $(e_1, ..., e_n)$ are regarded as applications of an n-ary constructor $Tuple_n$ to $n$ arguments $e_1, ..., e_n$. We use $v, w, x, y, z$ for denoting variables, $p, q, r$ for patterns, and $e, e'$ for expressions. This extended $\lambda$ calculus is straight-forwardly mapped into the (ordinary) $\lambda$ calculus (see [Pey87]). In particular, $let \ x = e' \ in \ e$ translates into $(\lambda x.e)e'$. The semantic soundness of $\beta$ reduction thus gives the first equivalence:

$$\text{let}/1 \quad let \ x = e' \ in \ e \quad = \quad e[e'/x]$$

This equivalence can be used to abstract common sub-expressions when used in a right-to-left direction, and to expand definitions in place when used in a left-to-right direction. The former operation will typically be useful at the end of the query transformation process, while the latter will be useful at its outset (in order to generate an overall expression to optimise).

Functions are defined by equations of the form $f = e$. If $f \in FV(e)$ i.e. if $f$ is recursively defined, the meaning of $f$ is given by the least fixed point of the higher-order (and non-recursive) function $\lambda f.e$ (see [Sch86]). This meaning may just be non-termination for some arguments e.g. for the function $f = \lambda x.not(f \ x)$. Thus, the semantic domain contains for each type $t$ an element $\bot_t$ which denotes 'no information' and represents a non-terminating computation (sometimes we omit the subscript $t$ when it can be inferred from context). For example, the boolean type consists of the elements, $True$, $False$ and $\bot_{Bool}$, where $\bot_{Bool}$ is less informative than both $True$ and $False$ (written $\bot_{Bool} \sqsubseteq True$ and $\bot_{Bool} \sqsubseteq False$) and where $True$ and $False$ are not information-wise comparable. The meaning of $f = \lambda x.not(f \ x)$ is then given by the least fixed point of the higher-order function $\lambda f.\lambda x.not(f \ x)$, and is just the function that maps all its arguments to $\bot_{Bool}$ i.e. $\lambda x.\bot_{Bool}$.

For the purposes of investigating query optimisation we use several items of information about expressions:

*Referential transparency.* This is a property enjoyed by our language and means that every occurrence of an expression denotes the same value in a given environment (an environment being a mapping of free variables to expressions).

*Termination.* The evaluation of an expression $e$ terminates if the value of $e$ contains no $\bot$ elements. In the sequel, whenever we say that an expression $e$ is *infinite* we mean that its value contains $\bot$; otherwise, we say that $e$ is *finite.* Determining whether $e$ is finite is of course undecidable in general. There is, however, a wide class of expressions whose evaluation is known to terminate, namely well-typed, non-recursive expressions: this is the *strong normalisation theorem* [Hin86]. Furthermore, it may be possible to construct a proof of the finiteness of an expression, perhaps using structural induction (see below) over its definition, and the user could be permitted to annotate the expression as such.

*Strictness of functions.* The order in which $\beta$ reduction is applied in $\lambda$ expressions is significant. Lazy evaluation (which we assume) ensures termination whenever possible by only evaluating the arguments to functions if needed by the function to return a result. A function is *strict* in an argument if that argument must be evaluated for the function to return a result. One way to characterise a strict function is to state that $f\bot = \bot$ (i.e. given a non-terminating argument, $f$ will not terminate either). Information about the strictness properties of a function can be derived from the known strictness properties of the built-in functions using *strictness analysis* [Cla85].

*Continuity of functions.* A function $f$ is said to be *continuous* if, for every sequence of values $d_1 \sqsubseteq d_2 \sqsubseteq ...$ in the domain of $f$, $f(\sqcup\{d_1, d_2, ...\}) = \sqcup\{f\,d_1, f\,d_2, ...\}$. In other words, continuous functions preserve least upper bounds. Any function defined in the $\lambda$ notation is continuous ([Sch86] Theorem 6.24). This has two important implications. First, it guarantees that any recursive definition has a unique meaning. Second, it means that when proving a proposition of the form $\forall x. f\,x = g\,x$, where $f$ and $g$ are continuous functions, induction over the structure of $x$ can be used to prove the equivalence even if $x$ is infinite i.e. contains $\bot$ elements. In the terminology of [Sch86] a proposition of the above form is an *inclusive predicate* for which fixpoint induction is valid. [Bir88] gives an accessible discussion of structural induction and uses it to prove equivalences over infinite lists and trees. We similarly use it to prove our equivalences over, possibly infinite, sets below.

Sets are important in our language, so we briefly recall their semantics (see [Pou93] for further details). The least element of the type consisting of sets of values of type $t$ is the set $\{\bot_t\}$. For example, the value

of $f = \lambda x.(fx) \cup (fx)$ is $\lambda x.\{\bot\}$, while the value of $nats = \lambda n.\{n\} \cup (nats(n+1))$ is $\lambda n.\{n, n+1, ...\bot\}$, so with $nats$ non-termination arises from the construction of an infinite set.

## 2.1 The type system

Our language is strongly, statically typed and supports a number of primitive types such as *Bool*, *Str* and *Num*. The user can declare new enumerated types and introduce new constants of such a type. For every user-defined enumerated type $T$, a built-in zero argument function $allT$ returns all constants of that type.

We will use as a running example a database that records results for the Winter Olympics. The user-defined enumerated types include *Comp* (competitor id), *Sex*, *Country*, *Category* (category of events) and *Event*, where:

$$
\begin{aligned}
allComp &= \{C1, C2, C3, ...\} \\
allSex &= \{Male, Female\} \\
allCountry &= \{Norway, Austria, France, ...\} \\
allCategory &= \{Alpine, Nordic, FigureSkating, ...\} \\
allEvent &= \{MensDownhill, WomensSlalom, ...\}
\end{aligned}
$$

Also supported are polymorphic product, list, set and function types. In particular, $(t_1, ..., t_n)$ is an n-product type for any types $t_1, ..., t_n$, $[t]$ is a list type and $\{t\}$ a set type for any type $t$, and $t_1 \to t_2$ the type of functions from a type $t_1$ to a type $t_2$. Note that the function type constructor $\to$ is right associative, so that $t_1 \to t_2 \to t_3$ and $t_1 \to (t_2 \to t_3)$ are synonymous.

We use the notation $e : t$ to indicate that an expression $e$ has type $t$. We also use letters from the start of the alphabet to indicate type variables in type expressions. For example, the infix 'compose' function, $\circ$, defined by $(f \circ g)x = f(g\,x)$ is of type $(b \to c) \to (a \to b) \to (a \to c)$, where the type variables $a, b$ and $c$ can be instantiated to any type.

The user can also declare sum types and introduce new constructors of such a type c.f. the list constructors $(:) : a \to [a] \to [a]$ and $[] : [a]$.

## 2.2 Built-in functions

The usual arithmetic $(+, -, *, /)$ and comparison $(==, != =, <, <=, >, >=)$ operators are built-in[1]. These operators may be written either infix or prefix (in which case they are bracketed e.g. $(+)\,1\,2$). These operators are of necessity strict in both their arguments. For example, the value of $e == e'$ will be $\bot$ if either $e$ or $e'$ has value $\bot$: operationally, both operands of $==$ must be evaluated in order to determine if they are equal, and if either yields 'no information' so does the overall evaluation of the equality test. The other main

---

[1]Note that $==$ is the syntactic equality operator as opposed to equality, $=$, in the semantic domain.

comparison operator, $<$, works in a similar fashion and hinges on an alphanumeric ordering of constants and constructors, with a left-to-right comparison of the arguments of the latter. Functions cannot be meaningfully compared, while sets are converted to lists (by the operation *set_to_list* below) for comparison.

The 3-argument conditional function *if* is also built-in and has the following semantics:

$$if \perp_{Bool} x \ y \ = \perp$$
$$if \ True \ x \ y \ = x$$
$$if \ False \ x \ y \ = y$$

Thus, *if* is strict in its first argument, but not in its second and third arguments. Logical operators can be defined in terms of *if* as follows:

$$and = \lambda x.\lambda y.if \ x \ y \ False$$
$$or \ = \lambda x.\lambda y.if \ x \ True \ y$$
$$not \ = \lambda x.if \ x \ False \ True$$

Consequently these too are strict in their first argument. We will also require on occasion counterparts to *and* and *or* that are commutative and '$\perp$-avoiding'; operationally, $\vee$ and $\wedge$ are implemented by evaluating their operands in parallel:

| | | | | |
|---|---|---|---|---|
| *True* $\vee$ *y* | $=$ *True* | | *False* $\wedge$ *y* | $=$ *False* |
| *x* $\vee$ *True* | $=$ *True* | | *x* $\wedge$ *False* | $=$ *False* |
| *False* $\vee$ *y* | $=$ *y* | | *True* $\wedge$ *y* | $=$ *y* |
| *x* $\vee$ *False* | $=$ *x* | | *x* $\wedge$ *True* | $=$ *x* |
| $\perp \vee \perp$ | $= \perp$ | | $\perp \wedge \perp$ | $= \perp$ |

Two set-building functions are also built-in, the singleton-forming operator and the union operator:

$$\{ \_ \} \quad : a \rightarrow \{a\}$$
$$\_ \cup \_ \ : \{a\} \rightarrow \{a\} \rightarrow \{a\}$$

A set can also be represented by enumerating its elements, where $\{e_1, ..., e_n\}$ equals $\{e_1\} \cup ... \cup \{e_n\}$.

Breazu-Tannen *et al.* [Bre91] propose a function, $\Phi$, for folding a binary operator *op* into a finite set, requiring *e* and *op* to form a commutative-idempotent monoid on the return type of *f* in order for this definition to have a unique meaning:

$$\Phi \ f \ op \ e \ \{\} \qquad = e$$
$$\Phi \ f \ op \ e \ \{x\} \qquad = f \ x$$
$$\Phi \ f \ op \ e \ (s1 \cup s2) = op \ (\Phi \ f \ op \ e \ s1) \ (\Phi \ f \ op \ e \ s2)$$

In our case of possibly infinite sets the second equation has to be modified, giving $\phi$ below:

$$\phi \ f \ op \ e \ \{\} \qquad = e$$
$$\phi \ f \ op \ e \ \{x\} \qquad = if \ (x = \perp) \perp \ (f \ x)$$
$$\phi \ f \ op \ e \ (s1 \cup s2) = op \ (\phi \ f \ op \ e \ s1) \ (\phi \ f \ op \ e \ s2)$$

In operational terms, $\phi$ keeps distributing *f* to the elements of an infinite set for ever; in semantic terms, the modified definition ensures that $\phi$ is continuous. Thus, as one might expect, we cannot use $\phi$ to devise terminating cardinality or summation functions for infinite sets. Others have defined similar functions to $\Phi$ e.g. the 'pump' operator of FAD [Ban87] and the 'hom' operator of Machiavelli [Oho89], and [Bre91] gives a comparison of these.

We can now use $\phi$ to define a membership operator over possibly infinite sets:

$$x \ in \ s \ = \ \phi \ ((==) \ x) \ (\vee) \ False \ s$$

Thus, *in* returns *True* if any comparison of *x* with an element of *s* returns *True*, *False* if all comparisons of *x* with elements of *s* return *False*, and $\perp$ otherwise. So *in* only returns *False* for finite sets: in operational terms, *in* keeps on searching an infinite set for a value that equals *x* until it finds one. Also, since *in* depends on the results of equality tests, it will return $\perp$ if applied to higher-order sets i.e. sets of functions.

We can also use $\phi$ to define a deterministic conversion function from sets to lists, where *list_union* is defined in terms of a duplicate-eliminating *sort* function and a list append operator ++:

$$set\_to\_list \ s \qquad = \phi \ (\lambda x.[x]) \ list\_union \ []$$
$$list\_union \ xs \ ys \qquad = sort \ (xs \ ++ \ ys)$$

Given *set_to_list* we can determine the cardinality of a set, the sum of its elements etc. Finally note that since *set_to_list* depends on *sort* it will return $\perp$ if applied to higher-order or infinite sets.

## 2.3 User-defined functions

These can be specified using one or more equations rather than a single $\lambda$ abstraction, and can use pattern-matching to deconstruct their arguments. For example, we can define a function *foldl* which is similar to $\phi$ but which works over lists:

$$foldl \ f \ op \ e \ [] \qquad = e$$
$$foldl \ f \ op \ e \ (x:xs) = op \ (f \ x) \ (foldl \ f \ op \ e \ xs)$$

This function can be used to convert a list to a set:

$$list\_to\_set \ xs \qquad = foldl \ (\lambda x.\{x\}) \ (\cup) \ \{\} \ xs$$

0-ary set-valued functions can be used to represent bulk data, the assumption being that such functions are updatable by the insertion and deletion of values of the appropriate type. For our Winter Olympics Database examples, we will use functions which define: the sponsors of each country; the name, sex and country of each competitor; the category of each event; the set of competitors registered for each event; and the list of medalists in rank order for each event:

| | | |
|---|---|---|
| *sponsors* | : | $\{(Country, Str)\}$ |
| *comps* | : | $\{(Comp, Str, Sex, Country)\}$ |
| *events* | : | $\{(Event, Category)\}$ |
| *registered* | : | $\{(Event, \{Comp\})\}$ |
| *results* | : | $\{(Event, [Comp])\}$ |

# 3  The Algebra

Our algebra consists of three built-in operators. These are the $\{\_\}$ and $\_\cup\_$ operators already introduced and an operator *setmap* which has the following semantics:

*setmap*       $: (a \rightarrow \{b\}) \rightarrow \{a\} \rightarrow \{b\}$
*setmap f* $\{\}$    $= \{\}$
*setmap f* $\{x\}$    $= if\ (x = \bot)\ \{\bot\}\ (f\ x)$
*setmap f (s1* $\cup$ *s2)* = *(setmap f s1)* $\cup$ *(setmap f s2)*

*setmap* thus distributes a function of type $(a \rightarrow \{b\})$ over a set of type $\{a\}$ and returns the union of the results. Notice that *setmap f* is just $\phi f\ (\cup)\ \{\}$. [Bre91] similarly gives a finite-set version of *setmap* defined in terms of $\Phi$.

Two functions that frequently appear in algebras are *map* : $(a \rightarrow b) \rightarrow \{a\} \rightarrow \{b\}$ and *filter* : $(a \rightarrow Bool) \rightarrow \{a\} \rightarrow \{a\}$ e.g. in [Clu92, Bee92]. These functions generalise relational projection and selection. Although they could be built-in for efficiency purposes, *map* and *filter* can be defined in terms of *setmap* as follows, where *singleton f x* = $\{f\ x\}$:

*filter f s*    $=$ *setmap* $(\lambda x.if\ (f\ x)\ \{x\}\ \{\})\ s$
*map f s*     $=$ *setmap (singleton f) s*

A further operation that can be expressed using *setmap* is the join of two relations according to a selection function $f$ and a projection function $g$, *join* : $((a, b) \rightarrow Bool) \rightarrow ((a, b) \rightarrow c) \rightarrow \{a\} \rightarrow \{b\} \rightarrow \{c\}$:

*join f g r s*    $=$ *setmap* $(\lambda x.setmap(\lambda y.$
                   *if (f (x,y)) {g (x,y)} {}) s) r*

*join* subsumes the various flavours of join and product operations found in relational databases. It can operate upon infinite sets of structured tuples. In particular, for any $a \in r$ and $b \in s$, the value of $g(a,b)$ in *(join f g r s)* is *True* provided that $f(a, b)$ is *True*, regardless of the finiteness or otherwise of $r$ and $s$.

Finally, *join* obeys the following commutativity property, provided $r$ and $s$ are both finite or both infinite, where $f^*\ (x,y) = f(y,x)$:

*join f g r s*    $=$ *join* $f^*\ g^*\ r\ s$

## 3.1  Other set-theoretic operators

Other set-theoretic operators can be defined in terms of the operators above, although these too could be built-in for efficiency purposes. We give definitions

for two of these operators, since they raise some interesting issues. Operators such as *nest*, *unnest* and *powerset* are also easily defined in our language.

Set difference can be defined using *filter* and *in*:

*_minus_*      $: \{a\} \rightarrow \{a\} \rightarrow \{a\}$
*s1 minus s2*   $=$ *filter* $(\lambda x.not(x$ *in s2)) s1*

Thus *minus* will terminate if both *s1* and *s2* are finite, or if *s1* is finite and is a subset of *s2*.

Intersection can also be defined using *filter* and *in*:

*_inter_*       $: \{a\} \rightarrow \{a\} \rightarrow \{a\}$
*s1 inter s2*   $=$ *filter* $(\lambda x.x$ *in s2) s1*

However, this definition is not in general commutative e.g. $\{3\}$ *inter* $\{3, \bot\} = \{3\}$ whereas $\{3, \bot\}$ *inter* $\{3\} = \{3, \bot\}$. Clearly it is desirable for intersection to be commutative for optimisation purposes. To achieve this we can use $\wedge$:

*s1 inter s2*   $=$   *filter* $(\lambda x.(x$ *in s1)* $\wedge$ *(x in s2))*
                      *(s1* $\cup$ *s2)*

This definition is both less efficient and has worse termination properties than the first e.g. $\{3\}$ *inter* $\{3, \bot\}$ now gives $\{3, \bot\}$, but is nevertheless the one we assume for optimisation purposes. If, however, both *s1* and *s2* are known to be finite then the original definition can safely be used in its place.

## 3.2  Equivalences

We now give some equivalences for the functions defined above. Some of these are generalisations of well-known equivalences for relational databases [Jar84, Ull89]. The main class of equivalences which do not have counterparts in our language are the commutative laws for joins and products. However, if records [Oho89] are used instead of tuples, these equivalences also apply, subject to the provisos stated for *join* above regarding the termination of its arguments.

The first set of equivalences, with their stated provisos, follow easily from the definitions of the logical operators in Section 2.2:

| | | |
|---|---|---|
| if/1 | *if e1 (if e2 e3 e4) e4* | $=$ *if (e1 and e2) e3 e4* |
| if/2 | *if e1 e3 (if e2 e3 e4)* | $=$ *if (e1 or e2) e3 e4* |
| if/3 | *if (not e1) e2 e3* | $=$ *if e1 e3 e2* |
| if/4 | *f (if e1 e2 e3)* | $=$ *if e1 (f e2) (f e3)* |
| | provided *f* is strict | |
| and/1 | *e1 and e2* | $=$ *e2 and e1* |
| | provided *e1* = $\bot$ iff *e2* = $\bot$ | |
| and/2 | *e1* $\wedge$ *e2* | $=$ *e2* $\wedge$ *e1* |
| or/1 | *e1 or e2* | $=$ *e2 or e1* |
| | provided *e1* = $\bot$ iff *e2* = $\bot$ | |
| or/2 | *e1* $\vee$ *e2* | $=$ *e2* $\vee$ *e1* |
| not/1 | *not (not e1)* | $=$ *e1* |
| not/2 | *not (e1 or e2)* | $=$ *(not e1) and (not e2)* |

The ∪ and *inter* operators obey the expected properties of commutativity and associativity (in operational terms, this means that the two branches of a ∪ must be evaluated in parallel), while *inter* and *minus* distribute over ∪:

∪/1  *s1* ∪ *s2*  = *s2* ∪ *s1*
∪/2  *s1* ∪ *(s2* ∪ *s3)*  = *(s1* ∪ *s2)* ∪ *s3*
∪/3  *s1* ∪ *s2*  = *s1*  if *s2* ⊆ *s1*
∩/1  *s1 inter s2*  = *s2 inter s1*
∩/2  *s1 inter (s2 inter s3)* = *(s1 inter s2) inter s3*
∩/3  *(s1 inter s3)* ∪ *(s2 inter s3)* =
     *(s1* ∪ *s2) inter s3*
−/1  *(s1 minus s3)* ∪ *(s2 minus s3)* =
     *(s1* ∪ *s2) minus s3*

However, the following equivalences only hold subject to the stated provisos:

∩/4  *s1 inter s2*  = *s2*
     if *s2* ⊆ *s1*, provided *s1* is finite
−/2  *s1 minus s2*  = *s1*
     if *s1* ∩ *s2* = {}, provided *s1* and *s2* are finite
−/3  *s1 minus s2*  = {}
     if *s1* ⊆ *s2*, provided *s1* is finite

To illustrate the proviso associated with ∩/4, consider the sets *s1* = {*1,2,⊥*} and *s2* = {*1*}. Then *s2* ⊆ *s1*, but *s1 inter s2* = {*1,⊥*} ≠ *s2*.

The set membership operator obeys the following properties, subject to the stated provisos:

in/1  *e1 in* {}  = *False*
in/2  *e1 in* {*e2*}  = *e1* == *e2*
in/3  *e in (s1* ∪ *s2)*  = *(e in s1)* ∨ *(e in s2)*
in/4  *e in (s1 inter s2)* = *(e in s1)* ∧ *(e in s2)*
      provided *e, s1, s2* are finite
in/5  *e in (s1 minus s2)*= *(e in s1)* ∧ *not(e in s2)*
      provided *e, s1, s2* are finite

∪/3, ∩/4 and −/3 also allow us to simplify the following expressions involving the built-in functions *allT*, provided s is finite:

all/1  *s* ∪ *allT*  = *allT* ∪ *s*  = *allT*
all/2  *s inter allT*  = *allT inter s*  = *s*
all/3  *s minus allT*  = {}

A number of optimisations apply to *setmap*, and hence also to operators defined in terms of *setmap* such as *filter* and *map*:

setmap/1 *setmap f (s1* ∪ *s2)* =
         *(setmap f s1)* ∪ *(setmap f s2)*
setmap/2 *setmap f (setmap g s)* =
         *setmap (λx.setmap f (g x)) s*
setmap/3 *setmap (λx.if (x in s1) e1 e2) s* =
         *(setmap (λx.e1) (s inter s1))* ∪

         *(setmap (λx.e2) (s minus s1))*
         provided *s* is finite
setmap/4 *setmap (λx.setmap (λy.e) s2) s1* =
         *setmap (λy.setmap (λx.e) s1) s2*
         provided *x* ∉ *FV(s2)*, *y* ∉ *FV(s1)*, and
         *s1* and *s2* are both finite or both infinite

setmap/1 states that *setmap* distributes over ∪. setmap/2 states that two successive applications of *setmap* can be compressed into one application with a second nested within it. setmap/3 states when application of a set membership test can be replaced by a set union. Finally setmap/4 states when the nesting of one *setmap* within another commutes. A consequence of setmap/4 is that any two *setmaps* which successively iterate over the *same* set can commute. This is especially important for the optimisation of iteration over recursively defined sets.

These equivalences are proved by structural induction over the set arguments. Since sets are constructed by successive unions of singleton sets and the empty set, structural induction over a set *s* has two base cases which must first be proved: *s* = {} and *s* = {*e*}. The induction hypothesis is then that the given proposition holds for sets *s'* and *s''*, from which it remains to show that it holds for *s* = *s'* ∪ *s''*.

The main optimisations for *map* are to combine successive applications into one. In particular map/2 below corresponds to combining cascades of projections:

map/1  *map f (map g s)*  = *map (f ∘ g) s*
map/2  *map (λq.r) (map (λp.q) s)* = *map (λp.r) s*
       provided *FV(q)* ⊆ *FV(p)*

For *filter*, filter/1 below is a generalised cascade of selections, filter/2 and filter/3 combine successive applications of *filter* and *setmap* into a single *setmap*, and filter/4 states that selection distributes over difference:

filter/1 *filter f (filter g s)*  =
         *filter (λx.(g x) and (f x)) s*
filter/2 *setmap f (filter g s)*  =
         *setmap (λx.if (g x) (f x) {}) s*
filter/3 *filter g (setmap f s)*  =
         *setmap (λx.filter g (f x)) s*
filter/4 *filter f (s1 minus s2)* =
         *(filter f s1) minus (filter f s2)*
         provided *s1* and *s2* are finite and
         *(f a)* is finite for finite *a*

Finally we have the expected equivalences regarding combining selection with join (join/1 below) and distributing selection over join (join/2):

join/1  *filter f1 (join f2 g r s)* =
        *join (λ(x,y).f2(x,y) and f1(g(x,y))) g r s*

join/2　*join (λ(x,y).(f1 x) and (f2 y)) g r s =*
*setmap (λx.if (f1 x) (setmap (λy.*
*if (f2 y) {g(x,y)} {}) s) {}) r*
*provided s finite and (f a) finite for finite a*

In summary, most of the expected equivalences for the logical and set operators hold. In some cases, however, we require functions to be strict, or we require a priori knowledge about the termination properties of expressions. The provisos associated with the equivalences arise from the semantics of the built-in functions. Clearly, built-in functions with different semantics, e.g. sequential ∨, ∧ and ∪, will give rise to different provisos.

### 3.3 Transformation Principles

Essentially the same principles apply to our language as to relational algebra expressions [Ull89], except that they need to be successively applied starting from the outermost level of an expression and moving through to expressions nested within aggregation functions:

(i) use filter/1 in a right to left direction, to split up complex filter conditions;

(ii) perform *filter* as early as possible by commuting it with other applications of *setmap* (setmap/4), eliminating set membership tests (setmap/3), and distributing *filter* over ∪ (setmap/1), *minus* (filter/4) and *join* (join/2);

(iii) perform *map* as early as possible by distributing it over ∪ (setmap/1);

(iv) combine cascades of *setmaps* of various kinds into a single *setmap* (setmap/2, map/1, map/2, filter/1-3, join/1);

(v) at any stage during the above steps, simplify set unions, intersections and differences whenever possible by using ∪/3, ∩/4, −/*, all/*, in/*.

The final step of the transformation process is to abstract common subexpressions using let/1 in a right-to-left direction.

Note that there is no general heuristic about which direction to apply setmap/1, since the size of the result returned by *setmap* cannot be predicted in general. Note also that further application of setmap/4 can be made using physical-level knowledge such as expected sizes of sets and availability of indexes. Finally, note that we could also derive an equivalence that moves *map* through *join* in the special case that the former is a projection and the latter a cartesian product, but would be quite contrived. In any case, such an equivalence would go into category (iii) above.

## 4　Higher-Level Constructs

The algebraic equivalences discussed above are very fine-grained and low level. We now examine two additional sets of equivalences at a higher conceptual modelling and querying level: those for set abstractions and those for functions with known inverses. Our reasons for doing so are two-fold. Firstly, both these constructs are commonly found database languages and we wish to extend optimisations identified by others (see Section 5) to our richer computational environment. Secondly, we conjecture that optimising first at this conceptual level is likely to be more efficient than proceeding directly to logical-level optimisations.

### 4.1　Set abstractions

The syntax of set abstractions is as follows:

| | | |
|---|---|---|
| *set_abstraction* | = | *"{" expr "∣" qualifiers "}"* |
| *qualifiers* | = | *qualifier ∣ qualifier ";" qualifiers* |
| *qualifier* | = | *generator ∣ filter* |
| *generator* | = | *pattern "∈" expr* |
| *filter* | = | *expr* |

For example, the following equations define a set *father* given a set *parent::{(Person,Person)}* and a set *mother::{(Person,Person)}* and a recursive set *anc*:

*father = {t ∣ t ∈ parent; not(t in mother)}*
*anc　= parent ∪ {(a,d) ∣ (a,d1) ∈ parent;*
*(a1,d) ∈ anc; a1 == d1}*

Optimisation of set abstractions is important since these provide a unifying query formalism for relational, functional, and deductive languages. For example, the head of a set abstraction corresponds to the SELECT clause of an SQL query, the generators are correspond to the FROM clause, and the filters to the WHERE clause. Also, Trinder [Tri89] gives a translation of the relational calculus into list (as opposed to set) abstractions, [Pat90] notes that DAPLEX queries are easily translated into set abstractions, and in previous papers e.g. [Pou93] we have observed the syntactic and semantic correspondence between set-valued functions such as *father* and *anc* and the analogous Datalog predicates. However, set abstractions are just syntactic sugar for nested applications of *setmap* and *if*. In the interests of simplicity, we give the translation scheme, *T* below, only for the case that the patterns in generators are simple variables: the interested reader can find the full translation scheme in [Pou93]. In the translation equations below *Q* denotes a sequence of zero or more qualifiers:

*T⌈{e∣}⌉　　= {T⌈e⌉}*
*T⌈{e₁∣e₂; Q}⌉　= if (T⌈e₂⌉) (T⌈{e₁∣Q}⌉) {}*
*T⌈{e₁∣x ∈ e₂; Q}⌉ = setmap (λx.T⌈{e₁∣Q}⌉) (T⌈e₂⌉)*

For example, *father* above translates into:

*setmap (λt.if (not (t in mother)) {t} {}) parent*

Three classes of equivalences can be identified for set abstractions, those for qualifier interchange, for qualifier introduction/elimination, and for moving qualifiers into nested set abstractions. These equivalences can be proved by translating expressions into the extended λ calculus and using structural induction.

The following equivalences for interchanging the qualifiers in set abstractions have well-known counterparts for list abstractions with bag equality [Tri89]:

abs/1 $\{e|\ ...;\ p1 \in s1;\ p2 \in s2;\ Q\} =$
$\quad\quad \{e|\ ...;\ p2 \in s2;\ p1 \in s1;\ Q\}$
$\quad\quad$ provided $FV(p1) \cap FV(s2) =$
$\quad\quad\quad\quad FV(p2) \cap FV(s1) = \{\}$
abs/2 $\{e|\ ...;\ p \in s;\ f;\ Q\}\quad = \{e|\ ...;\ f;\ p \in s;\ Q\}$
$\quad\quad$ provided $FV(p) \cap FV(f) = \{\}$
abs/3 $\{e|\ ...;\ f;\ g;\ Q\}\quad\quad = \{e|\ ...;\ g;\ f;\ Q\}$

abs/1 states that generators can be interchanged. It follows directly from setmap/4 and has the same proviso that s1 and s2 are both finite or both infinite. abs/2 states that a generator and a filter can be interchanged, and it requires both that *s* is finite and that *f* terminates, otherwise non-termination may be introduced. Of course, if we do not mind improving the termination properties of an expression, the rule may be used in a right-to-left direction if *s* is known to be finite, and in a left-to-right direction if *f* is known to terminate. abs/3 states that two filters can be interchanged. Its proof requires if/1 and and/1, and consequently this equivalence holds only if *f* fails to terminate whenever *g* does.

Numerous equivalences can be identified for eliminating qualifiers, of which the following is a representative sample:

abs/4 $\{e|\ ...;\ f;\ g;\ Q\}\quad\quad = \{e|\ ...;\ f\ and\ g;\ Q\}$
abs/5 $\{e|...;\ x \in \{e'\};\ Q\} = \{e[e'/x]|\ ...;\ Q[e'/x]\}$
$\quad\quad$ provided $x \notin BV(Q)$
abs/6 $\{e|\ ...;\ p \in s1;\ p\ in\ s2;\ Q\} =$
$\quad\quad\quad \{e|\ ...;\ p \in (s1\ inter\ s2);\ Q\}$

abs/4 states that two filters can be compressed into one: its proof follows directly from if/1. abs/5 states that a generator over a singleton can be eliminated: its proof follows from the semantic soundness of β reduction. abs/6 states that a filter can be eliminated: its proof follows from the definition of *inter* and only holds if both *s1* and *s2* are finite.

The third set of equivalences governs the moving of qualifiers into and out of nested set abstractions:

abs/7 $\{e|\ \ ...;\ p \in s;\ Q\} = \{e|\ ...;\ p \in \{p \mid p \in s\};\ Q\}$
abs/8 $\{e|...;\ p \in \{p|Q\};\ f;\ Q'\} =$

$\quad\quad \{e|...;\ p \in \{p|Q;\ f\};\ Q'\}$
$\quad\quad$ provided $FV(f) \subseteq FV(p)$

More sophisticated forms of these are

$\{e|\ ...;\ p \in s;\ Q\} = \{e|\ ...;\ p \in \{p' \mid p' \in s\};\ Q\}$
$\{e|\ ...;\ p \in \{p'|Q\};\ f;\ Q'\} =$
$\quad\quad \{e|\ ...;\ p \in \{p'|Q;\ f'\};\ Q'\}$

where $p'$ is obtained from $p$ by a renaming of variables and $f'$ is obtained from $f$ by the same renaming.

## 4.2 Transformation of set abstractions

The following transformation principles can be applied successively, moving inwards from the outermost set abstraction to nested set abstractions c.f. the lower-level transformations of 3.3:

(i) use abs/4 in a right to left direction to split up complex filter conditions;

(ii) perform filters as early as possible by interchanging them with generators and other filters using abs/2 and abs/3;

(iii) interchange groups consisting of a generator and its dependent filters by using abs/1-3, according to the expected efficiency of evaluating the group (based on physical-level knowledge such as expected sizes of sets and availability of indexes);

(iv) eliminate redundant qualifiers using abs/4-6;

(v) at any stage during the above steps, simplify set unions, intersections and differences whenever possible by using the equivalences of Section 3;

(vi) pass filters into preceding, nested, set abstractions using abs/8.

Finally, abstract common subexpressions using let/1 in a right-to-left direction.

We illustrate these principles via two queries. The first requires the countries of women competitors who won alpine events. A naive formulation iterates through all countries, events, results and competitors, and then specifies the join condition:

$\{c \mid c \in allCountry;\ (ev,cat) \in events;$
$\quad (ev',nums) \in results;$
$\quad (num,name,sex,compc) \in comps;$
$\quad ev' == ev\ and\ cat == Alpine\ and\ num ==$
$\quad head(nums)\ and\ sex == Female\ and\ compc == c\}$

Applying abs/4 in a right to left direction, followed by a promotion of filters gives:

$\{c \mid c \in allCountry;\ (ev,cat) \in events;\ cat == Alpine;$
$\quad (ev',nums) \in results;\ ev' == ev;$
$\quad (num,name,sex,compc) \in comps;$
$\quad num == head(nums);\ sex == Female;\ compc == c\}$

Interchange of groups of generators and their dependent filters gives:

$\{c \mid (ev,cat) \in events; cat == Alpine;$
    $(ev',nums) \in results; ev' == ev;$
    $(num,name,sex,compc) \in comps;$
    $num == head(nums); sex == Female;$
    $c \in allCountry; compc == c\}$

or, alternatively:

$\{c \mid (num,name,sex,compc) \in comps;$
    $sex == Female; c \in allCountry; compc == c;$
    $(ev',nums) \in results; num == head(nums);$
    $(ev,cat) \in events; cat == Alpine; ev' == ev\}$

Compressing filters, and removing $c \in allCountry;$ $compc == c$ by using in/2, followed by abs/6 and all/2, gives the following for the first of these alternatives:

$\{c \mid (ev,cat) \in events; cat == Alpine;$
    $(ev',nums) \in results; ev' == ev;$
    $(num,name,sex,compc) \in comps;$
    $num == head(nums)$ and $sex == Female;$
    $c \in \{compc\}\}$

Finally, using abs/5 gives:

$\{compc \mid (ev,cat) \in events; cat == Alpine;$
    $(ev',nums) \in results; ev' == ev;$
    $(num,name,sex,compc) \in comps;$
    $num == head(nums)$ and $sex == Female\}$

A similar process for the second alternative gives:

$\{compc \mid (num,name,sex,compc) \in comps;$
    $sex == Female;$
    $(ev',nums) \in results; num == head(nums);$
    $(ev,cat) \in events; cat == Alpine$ and $ev' == ev\}$

The second query requires the competitors sponsored by Atomic who won alpine events. A naive formulation iterates through all events, all results and all tuples of a join of competitors with sponsors over country, and then specifies a further join condition:

$\{c \mid (ev,cat) \in events; (ev',nums) \in results;$
    $(c,spname) \in$
        $\{(c,spname) \mid (c,name,sex,compc) \in comps;$
        $(spname,spc) \in sponsors; compc == spc\}$
    $spname == \text{``Atomic''}$ and $ev' == ev$ and
    $cat == Alpine$ and $c == head(nums)\}$

Applying abs/4 in a right to left direction followed by filter promotion gives two alternatives, one being:

$\{c \mid (ev,cat) \in events; cat == Alpine;$
    $(ev',nums) \in results; ev' == ev;$
    $(c,spname) \in$
        $\{(c,spname) \mid (c,name,sex,compc) \in comps;$
        $(spname,spc) \in sponsors; compc == spc\}$
    $spname == \text{``Atomic''}; c == head(nums)\}$

Passing the last two filter conditions into the preceding set abstraction using abs/8 gives:

$\{c \mid (ev,cat) \in events; cat == Alpine;$
    $(ev',nums) \in results; ev' == ev;$
    $(c,spname) \in$
        $\{(c,spname) \mid (c,name,sex,compc) \in comps;$
        $(spname,spc) \in sponsors; compc == spc;$
        $spname == \text{``Atomic''}; c == head(nums)\}\}$

Then performing filter promotion within the nested set abstraction followed by filter compression gives:

$\{c \mid (ev,cat) \in events; cat == Alpine;$
    $(ev',nums) \in results; ev' == ev;$
    $(c,spname) \in$
        $\{(c,spname) \mid (c,name,sex,compc) \in comps;$
        $c == head(nums);$
        $(spname,spc) \in sponsors;$
        $compc == spc$ and $spname == \text{``Atomic''}\}\}$

or

$\{c \mid (ev,cat) \in events; cat == Alpine;$
    $(ev',nums) \in results; ev' == ev;$
    $(c,spname) \in$
        $\{(c,spname) \mid (spname,spc) \in sponsors;$
        $spname == \text{``Atomic''};$
        $(c,name,sex,compc) \in comps;$
        $compc == spc$ and $c == head(nums)\}\}$

The second alternative is optimised similarly.

### 4.3 Functions with known inverses

In this section we consider how use may be made of information about the inverse relationship between pairs of functions. Such information may be readily available in languages which use a functional or object-oriented data model, and may also be available from other sources e.g. via proofs supplied by the programmer [Har92]. In the case of extensional functions, inverses typically correspond to fast access paths provided by indexes. An implicit assumption for the equivalences we give below is that the functions and their inverses are strict, which is necessary in order for the equivalences to hold.

Given two functions $f : s \to t$ and $f^{-1} : t \to s$ such that $f^{-1}(f\ e) = e$ for all expressions $e : s$ and $f(f^{-1}\ e) = e$ for all expressions $e : t$, then:

inv/1  $(f\ e1) == e2$  $=$  $e1 == (f^{-1}\ e2)$
inv/2  $(f\ e1)\ in\ s$  $=$  $e1\ in\ (map\ f^{-1}\ s)$

Given a function $f : s \to t$ which is many-to-one and a function $f^{-1} : t \to \{s\}$ such that for all $e : s$, $e \in$ $f^{-1}(f\ e)$ and for all $e : t$, $\forall u \in (f^{-1}e).f\ u = e$, then:

inv/3  $(f\ e1) == e2$  $=$  $e1\ in\ (f^{-1}\ e2)$
inv/4  $(f\ e1)\ in\ s$  $=$  $e1\ in\ (setmap\ f^{-1}\ s)$

Clearly inv/3 and inv/4 are also applicable given a one-to-many function $f : s \rightarrow \{t\}$ and its inverse $f^{-1} : t \rightarrow s$, since in this case the roles of $f$ and $f^{-1}$ above are reversed. Finally, given a many-to-many function $f : s \rightarrow \{t\}$ and a function $f^{-1} : t \rightarrow \{s\}$ such that for all $e : s$, $\forall u \in (f\,e).e \in (f^{-1}u)$ and for all $e : t$, $\forall u \in (f^{-1}e).e \in (f\,u)$, then:

inv/5   *e2 in (f e1)     = e1 in (f⁻¹ e2)*

### 4.4 Transformation of functions with inverses

The following transformation principles should be applied at successive levels of nesting, and common sub-expressions should be abstracted out at the end of the process:

(i) use abs/4 in a right to left direction to break up complex filter conditions;

(ii) use inv/1-5 to generate tests for set-membership and equality on variables i.e. tests of the form $v == e$ and $v\,in\,e$ where $v$ is a variable;

(iii) apply transformation principles (ii) - (v) for set abstractions from Section 4.2.

We again illustrate these principles on our Winter Olympics database, assuming the following functions and inverses which can be defined in terms of the base functions of Section 2.3:

*country_of* : *Comp → Country*
*team_of* : *Country → {Comp}  //= country_of⁻¹*
*competes_in* : *Comp → {Event}*
*comps_in* : *Event → {Comp}  //= competes_in⁻¹*
*winner* : *Event → Comp*
*sex* : *Comp → Sex*
*sex⁻¹* : *Sex → {Comp}*
*category* : *Event → Category*
*events* : *Category → {Event}  //= category⁻¹*

The query is the same as the first query of Section 4.2, i.e. 'which countries have female Alpine gold medalists?', and can be expressed as follows:

$\{c \mid c \in$ *allCountry; comp $\in$ allComp; ev $\in$ allEvents;*
  *Alpine == (category ev) and Female == (sex comp)*
  *and comp == (winner ev)*
  *and c == (country_of comp)*$\}$

Phase (i), which repeatedly applies abs/4 to break up the filter condition, gives:

$\{c \mid c \in$ *allCountry; comp $\in$ allComp; ev $\in$ allEvents;*
  *Alpine == (category ev); Female == (sex comp);*
  *comp == (winner ev); c == (country_of comp)*$\}$

During phase (ii) inv/1-5 are applied, and we reach the point:

$\{c \mid c \in$ *allCountry; comp $\in$ allComp; ev $\in$ allEvents;*
  *ev in (events Alpine); comp in (sex⁻¹ Female);*
  *comp == (winner ev); c == (country_of comp)*$\}$

at which we have a choice: whether or not to apply inv/1 to $c == (country\_of\ comp)$. If we do not do so, we move on to phase (iii), where the filters are removed using abs/6:

$\{c \mid ev \in$ *allEvents inter (events Alpine);*
  *comp $\in$ allComp inter (sex⁻¹ Female)*
    *inter {winner ev};*
  *c $\in$ allCountry inter {country_of comp}*$\}$

Finally, the lower level optimisations (in particular, all/2) reduce the size of the sets over which we iterate, giving us our first query plan:

$\{c \mid ev \in$ *events Alpine;*
  *comp $\in$ (sex⁻¹ Female) inter {winner ev};*
  *c $\in$ {country_of comp}*$\}$

Alternatively, applying inv/1 again gives:

$\{c \mid c \in$ *allCountry; comp $\in$ allComp; ev $\in$ allEvents;*
  *ev in (events Alpine); comp in (sex⁻¹ Female);*
  *comp == (winner ev); comp in (team_of c)*$\}$

and removing the filters using abs/6 gives:

$\{c \mid c \in$ *allCountry;*
  *ev $\in$ allEvents inter (events Alpine);*
  *comp $\in$ allComp inter (sex⁻¹ Female) inter*
    *{winner ev} inter (team_of c)*$\}$

Again, all/2 reduces the size of the sets over which we iterate, giving a second query plan:

$\{c \mid c \in$ *allCountry; ev $\in$ (events Alpine);*
  *comp $\in$ (sex⁻¹ Female) inter {winner ev}*
    *inter (team_of c)*$\}$

These two query plans differ only in that the extra use of inv/1 during the construction of the second query plan did not allow the iteration through *allCountry* to be eliminated. Thus, it is likely that the first plan would be chosen for execution after applying some physical-level heuristics. We finally observe that both query plans correspond to the first plan of Section 4.2 for the same query. The second plan of Section 4.2 is not generated here due to the non-availability of an inverse for the *winner* function.

## 5  Related Work

Our framework draws considerably from Breazu-Tannen *et al.* [Bre91] who proposed a programming paradigm based upon structural recursion over sets. A major motivating factor behind this paradigm is that

it facilitates the optimisation of database programs. Indeed, several optimisations are stated, including filter/2, filter/3 and join/2 above, while structural induction is used as a proof technique. However, this research is in the context of terminating functions over finite sets and a single level at which equivalences can be expressed. In contrast, we have richer semantic domains which include $\perp$, and can thus address termination issues; we also allow for the possibility that sets may be infinite; and we have investigated equivalences expressed at several levels of abstraction.

Cluet and Delobel [Clu92] propose a query optimisation formalism for $O_2$ based upon classes and algebraic query rewriting. One assumption made is that methods have no side-effects, although it would seem difficult in practice to guarantee that this condition holds. A select-project-join algebra is discussed, and the introduction of types into algebraic expressions facilitates its optimisation by allowing the introduction of functions that enumerate the constants of a type. Our use of the functions $allT$ and the equivalences inv/1-5 is analogous. The approach of [Clu92] retains a separation of DML and programming language; thus, the possibility of non-termination or infinite data structures are not considered.

Several other groups have introduced object algebras and strategies for their optimisation [Dem94, Lie92, Sha89, Sto91, Van91] with analogous equivalences to those we propose here. In general these algebras are either computationally incomplete or support optimisations for only a subset of their operators. Also, some provide only limited facilities for optimising user-defined data types; while others allow few algebraic transformations to be applied to an expression without changing its value.

Finally, the optimisation of functional database languages has been examined by several other researchers e.g. [Tri89, Bee90, Erw91, Pat90, Hey91]. Trinder [Tri89] advocates analogues of abs/1-3 for list (as opposed to set) abstractions. These equivalences are justified by assuming bag equality over lists. Implicit assumptions made are that lists are finite (otherwise, the equivalent of abs/1 would not hold for example) and that functions over their elements are terminating (otherwise abs/3 would not hold). Trinder also proposes *filter hiding*, which corresponds to a combination of our abs/7 and abs/8. [Hey91] too is concerned with the optimisation of list-valued expressions, and in this context proposes combining unary expressions (analogous to our map/1 and filter/1-3), eliminating iteration over functions of the form $allT$, and using information regarding indexes to select preferred query paths. [Pat90] discusses the optimisation of DAPLEX queries, essentially using abs/1-3 and inv/1-5. [Bee90, Erw91] discuss the optimisation of FP-like functional database languages, also highlighting the suitability of using functional language for DBPL optimisation. In particular, [Erw91] develops equivalences over map, filter and a $\Phi$-like aggregation operator in the context of strict functions over finite sets; a set of inverse equivalences are also given, including inv/1, and others that we have not discussed here.

# 6 Conclusions

We have investigated algebraic query optimisation techniques in the context of a functional DBPL furnished with a set bulk data type. We have examined the extent to which prior work on the optimisation of relational languages can be utilised. The declarative nature of our language has enabled us to avoid the problems associated with side effects, whilst its well defined semantics provides a framework in which to show formally termination properties of expressions and equivalences between expressions. We have identified caveats to several well-known equivalences in this richer computational paradigm. For processing tasks such as aggregation and transitive closure, our optimisations can be fully exploited for all sub-expressions of a query since there is no dichotomy between the optimisation of 'programs' and 'DML statements'.

Although developed in the context of a functional language, our findings are directly applicable to other DBPLs operating over sets. Conversely, we can incorporate equivalences discovered by others into our formalism. Finally, although we have concentrated upon showing equivalences relevant to the set data type, the same approach can be used for other, possibly user-defined, data types: see for example the equivalences shown in [Bir88] for list and tree data types.

We have not yet considered heuristics to reduce the query search space. Also, in a computationally complete language the user can easily 'program around' predefined functions and thus thwart the optimisation process. One ongoing issue is the optimisation of retrieval from recursive sets. We discussed this in [Pou93] where we addressed the optimisation of a class of set-valued functions called *selectors*. Selectors generalise the inverse functions of a functional data model by allowing associative look-up into n-ary, as opposed to just binary, relations. We proposed a magic-sets like rewriting of the definitions of intentional selectors, given specific look-up patterns: we are now implementing these ideas. Finally, it is clear that the user must be provided with sophisticated tools if they are to aid the optimisation process. Such tools have already been developed for functional languages, examples being strictness analysis [Cla85] and Cambridge LCF [Pau87], which can be used to prove properties of functions such as equivalence and termination.

## Acknowledgements

## References

[Alb91]  Albert, J. *Algebraic properties of bag data types*, Proc. 17th VLDB Conference, 1991.

[Ban87]  Bancilhon, F. et al. *FAD, a powerful and simple database language*, Proc. 13th VLDB Conference, 1991.

[Bee90]  Beeri, C. and Kornatzky, Y. *Algebraic optimization of object-oriented query languages*, Proc. 3rd ICDT, LNCS 470, Springer-Verlag, 1990.

[Bee92]  Beeri, C. and Milo, T. *Functional and predicative programming in OODBs*, Proc. ACM PODS, 1992.

[Bir88]  Bird, R. and Wadler, P. *Introduction to functional programming*, Prentice-Hall, 1988.

[Bre91]  Breazu-Tannen, V. Buneman, P. and Naqvi, S. *Structural recursion as a query language*, Proc. 3rd DBPL, Morgan-Kaufman, 1991.

[Cla85]  Clack, C. and Peyton-Jones, S. *Strictness analysis - a practical approach*, Proc. FPCA, Springer-Verlag LNCS 201, 1985.

[Clu92]  Cluet, S. and Delobel, C. *A General Framework for the Optimization of Object-Oriented Queries*, Proc. ACM SIGMOD, 1992.

[Dem94]  Demuth, B. Geppert, A. Gorchs, T. *Algebraic query optimisation in the CoOMS structurally object-oriented database system*, in Query Processing for Advanced Database Systems, eds. Freytag, Maier and Vossen, Morgan Kaufman, 1994.

[Erw91]  Erwig, M. and Lipeck, U. *A functional DBPL revealing high level optimizations*, Proc. 3rd DBPL, Morgan-Kaufman, 1991.

[Har92]  Harrison, P. and Khoshafian, S. *The mechanical transformation of data types*, The Computer Journal, 35(2), 1992.

[Hey91]  Heytens, M.L. and Nikhil, R.S. *List comprehensions in Agna, a parallel persistent object system*, Proc. FPCA, Springer-Verlag LNCS 523, 1991.

[Hin86]  Hindley, J.R. and Seldin, J.P. *Introduction to combinators and λ-calculus*, C.U.P., 1986.

[Jar84]  Jarke, M. and Koch, J. *Query optimisation in database systems*, ACM Computing Surveys, 16(2), 1984.

[Lie92]  Lieuwin, D. and Dewitt, D. *A transformation - based approach to optimizing loops in database programming languages*, Proc. ACM SIGMOD, 1992.

[Oho89]  Ohori, A. Buneman, P. and Breazu-Tannen, V. *Database programming in Machiavelli - a polymorphic language with static type inference*, Proc. ACM SIGMOD, 1989.

[Pat90]  Paton, N.W. and Gray, P.M.D. *Optimising and executing DAPLEX queries using Prolog*, The Computer Journal, 33(6), 1990.

[Pau87]  Paulson, L.C. *Logic and computation: interactive proof with Cambridge LCF*, Cambridge University Press, 1986.

[Pey87]  Peyton-Jones, S.L. *The Implementation of Functional Programming Languages*, Prentice Hall, 1987

[Pou93]  Poulovassilis, A. and Small, C. *A domain-theoretic approach to integrating logic and functional database languages*, Proc. 19th VLDB Conference, 1993.

[Sch86]  Schmidt, D.A. *Denotational Semantics*, Allyn and Bacon, 1986.

[Sha89]  Shaw, G.B. and Zdonik, S.B. *An object oriented query algebra*, Proc. 2nd DBPL, Morgan-Kaufman, 1989.

[Sto91]  Stonebraker, M. *Managing persistent objects in a multi-level store*, Proc. ACM SIGMOD, 1991.

[Tri89]  Trinder, P. *A functional database*, D. Phil. Thesis, Oxford University, 1989.

[Ull89]  Ullman, J.D. *Principles of database and knowledge-base systems*, Vol. 2, Computer Science Press, 1989.

[Van91]  Vandenburg, S.L. and DeWitt, D.J. *Algebraic support for complex objects with arrays, identity and inheritance*, Proc. ACM SIGMOD, 1991.