

Efficient Optimization of Iterative Queries

Leonidas Fegaras
Oregon Graduate Institute

A New Query Algebra Based on Folds

Folds

- can be defined for a large number of bulk data types;
- can capture most algebraic operators, such as *map*, *filter*, and *join*.

Programming with folds supports:

- calculation-based program optimizations;
- deforestation and loop fusion;
- equational reasoning and theorem proving;
- program synthesis.

So What is a Fold?

Algebraic Types

A **fold** is the natural control structure for a freely constructed algebraic type (a sums-of-products type):

$$T(\alpha_1, \dots, \alpha_p) = \begin{array}{l} C_1(t_{1,1}, \dots, t_{1,m_1}) \\ | \dots \\ C_n(t_{n,1}, \dots, t_{n,m_n}) \end{array}$$

where $\alpha_1, \dots, \alpha_p$ are type variables;

the C_i are value constructors;

and $t_{i,j}$ are either type variables

or instantiations of algebraic types

or the type $T(\alpha_1, \dots, \alpha_p)$ itself.

Examples of Algebraic Types

boolean = False | True

prod(α , β) = Pair(α , β)

list(α) = Nil | Cons(α , list(α))

nat = Zero | Succ(nat)

tree(α) = Tip(α) | Node(tree(α), tree(α))

bush(α) = Leaf(α) | Branch(list(bush(α)))

term(α) = Var(α) | Lambda(α , term(α)) | Apply(term(α), term(α))

person = Make_person(string, nat, string)
where string = list(nat)

The Fold Operator

- **Lists:**

If $g : \text{list}(\alpha) \rightarrow \beta$ can be expressed as:

$$\begin{aligned}g(\text{Nil}) &= f_n() \\g(\text{Cons}(a, l)) &= f_c(a, g(l))\end{aligned}$$

then g is a list fold: $g = \text{fold}^{\text{list}}(f_n, f_c)$.

$$\begin{aligned}\text{fold}^{\text{list}}(f_n, f_c) \text{Nil} &= f_n() \\ \text{fold}^{\text{list}}(f_n, f_c) \text{Cons}(a, l) &= f_c(a, \text{fold}^{\text{list}}(f_n, f_c) l)\end{aligned}$$

- **Natural numbers:**

$$\begin{aligned}\text{fold}^{\text{nat}}(f_z, f_s) \text{Zero} &= f_z() \\ \text{fold}^{\text{nat}}(f_z, f_s) \text{Succ}(i) &= f_s(\text{fold}^{\text{nat}}(f_z, f_s) i)\end{aligned}$$

- **Booleans:**

$$\begin{aligned}\text{fold}^{\text{boolean}}(f_f, f_t) \text{False} &= f_f() \\ \text{fold}^{\text{boolean}}(f_f, f_t) \text{True} &= f_t()\end{aligned}$$

Or $\text{fold}^{\text{boolean}}(f_f, f_t) x \equiv \mathbf{\text{if } x \text{ then } f_t() \text{ else } f_f()}$.

- **Trees:**

$$\begin{aligned}\text{fold}^{\text{tree}}(f_t, f_n) \text{Tip}(a) &= f_t(a) \\ \text{fold}^{\text{tree}}(f_t, f_n) \text{Node}(l, r) &= f_n(\text{fold}^{\text{tree}}(f_t, f_n) l, \text{fold}^{\text{tree}}(f_t, f_n) r)\end{aligned}$$

Example: List Append

$$\begin{aligned}\text{append}(\text{Nil}, y) &= y \\ \text{append}(\text{Cons}(a, l), y) &= \text{Cons}(a, \text{append}(l, y))\end{aligned}$$

Append as a fold:

$$\text{append}(x, y) = \text{fold}^{list}(\lambda().y, \lambda(a, r).\text{Cons}(a, r)) x$$

Variable r is an *accumulative result variable*.

Other Examples:

$$\text{length}(x) = \text{fold}^{list}(\lambda().\text{Zero}, \lambda(a, r).\text{Succ}(r)) x$$

$$\text{reverse}(x) = \text{fold}^{list}(\lambda().\text{Nil}, \lambda(a, r).\text{append}(r, \text{Cons}(a, \text{Nil}))) x$$

$$x + y = \text{fold}^{nat}(\lambda().y, \lambda(r).\text{Succ}(r)) x$$

$$x \times y = \text{fold}^{nat}(\lambda().\text{Zero}, \lambda(r).y + r) x$$

$$\text{even}(x) = \text{fold}^{nat}(\lambda().\text{False}, \lambda(r).\text{not}(r)) x$$

$$x \wedge y = \text{fold}^{boolean}(\lambda().\text{False}, \lambda().y) x$$

$$\text{flat}(x) = \text{fold}^{tree}(\lambda(i).\text{Cons}(i, \text{Nil}), \lambda(l, r).\text{append}(l, r)) x$$

$$p.\text{address} = \text{fold}^{person}(\lambda(n, s, a).a) p$$

The Fold Operator for any Algebraic Type T

For each constructor C_i of T we associate a *functor* E_i^T :

e.g. $\text{Cons} : (\alpha \times \text{list}(\alpha)) \rightarrow \text{list}(\alpha)$

Then

$$E_{\text{cons}}^{\text{list}(\alpha)}(S) = \alpha \times S$$

$$E_{\text{cons}}^{\text{list}(\alpha)}(g) = \text{id} \times g = \lambda(a, r).(a, g(r))$$

The **fold** operator for any algebraic type T is:

$$\text{fold}^T(\bar{f}) \circ C_i = f_i \circ E_i^T(\text{fold}^T(\bar{f}))$$

$$\begin{array}{ccc} T & \xleftarrow{C_i} & E_i^T(T) \\ \text{fold}^T(\bar{f}) \downarrow & & \downarrow E_i^T(\text{fold}^T(\bar{f})) \\ \beta & \xleftarrow{f_i} & E_i^T(\beta) \end{array}$$

Promotion Theorems

- **Lists:**

$$\begin{array}{l} \phi_n() = g(f_n()) \\ \phi_c(a, g(r)) = g(f_c(a, r)) \end{array} \hrule g(\text{fold}^{list}(f_n, f_c) x) = \text{fold}^{list}(\phi_n, \phi_c) x$$

- **Natural numbers:**

$$\begin{array}{l} \phi_z() = g(f_z()) \\ \phi_s(g(r)) = g(f_s(r)) \end{array} \hrule g(\text{fold}^{nat}(f_z, f_s) x) = \text{fold}^{nat}(\phi_z, \phi_s) x$$

- **Booleans:**

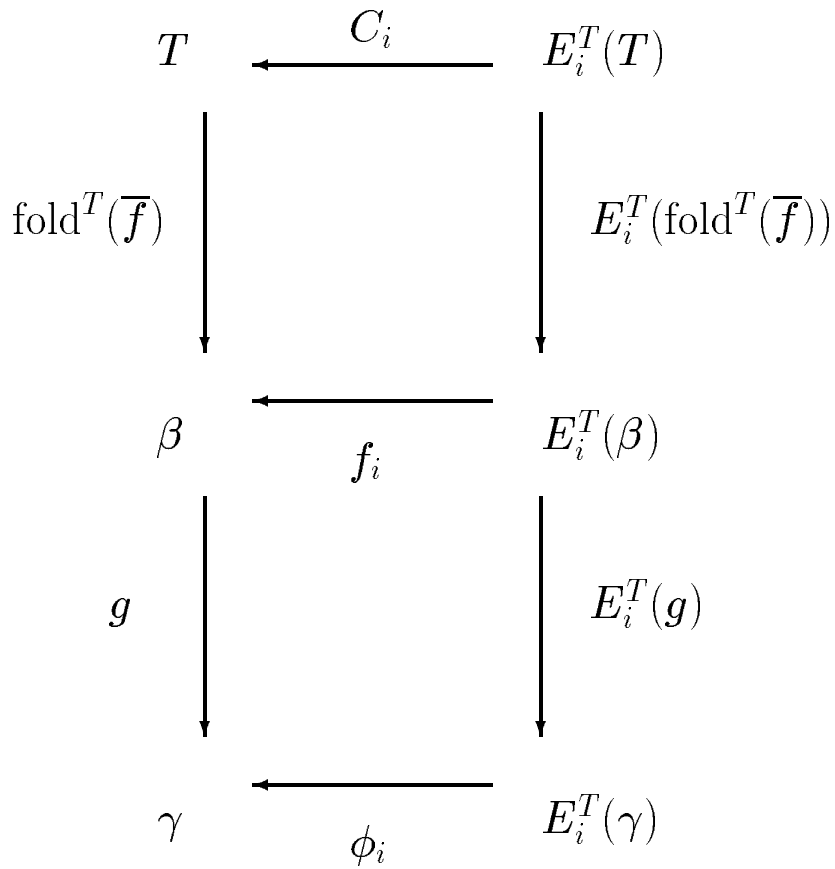
$$\begin{array}{l} \phi_f() = g(f_f()) \\ \phi_t() = g(f_t()) \end{array} \hrule g(\text{fold}^{boolean}(f_f, f_t) x) = \text{fold}^{boolean}(\phi_f, \phi_t) x$$

Or equivalently:

$$g(\text{if } x \text{ then } f_t() \text{ else } f_f()) = \text{if } x \text{ then } g(f_t()) \text{ else } g(f_f())$$

The Promotion Theorem for any Algebraic Type T

$$\frac{\forall i : \phi_i \circ E_i^T(g) = g \circ f_i}{g \circ \text{fold}^T(\bar{f}) = \text{fold}^T(\bar{\phi})}$$



Uniqueness Property

$$\forall i : g \circ C_i = \phi_i \circ E_i^T(g) \iff g = \text{fold}^T(\bar{\phi})$$

Sets

Set constructors: Emptyset and Insert

Set selector: split

Set fold:

$$\text{fold}^{set}(f_e, f_s) s = \begin{cases} \text{if } s = \text{Emptyset} \text{ then } f_e() \\ \text{else let } (a, r) = \text{split}(s) \text{ in } f_s(a, \text{fold}^{set}(f_e, f_s) r) \end{cases}$$

For example:

$$\mathbf{union}(x, y) = \text{fold}^{set}(\lambda().y, \lambda(a, r).\text{Insert}(a, r)) x$$

$$\mathbf{member}(e, x) = \text{fold}^{set}(\lambda().\text{False}, \lambda(a, r).(a = e) \mathbf{or} r) x$$

$$\begin{aligned} \mathbf{restrict}(s, f) = \\ \text{fold}^{set}(\lambda().\text{Emptyset}, \\ \lambda(a, r).\mathbf{if} f(a) \mathbf{then} \text{Insert}(a, r) \mathbf{else} r) s \end{aligned}$$

$$\begin{aligned} \mathbf{join}(x, y, M, J) = \\ \text{fold}^{set}(\lambda().\text{Emptyset}, \\ \lambda(a, r).\text{fold}^{set}(\lambda().r, \lambda(b, s).\mathbf{if} M(a, b) \mathbf{then} \text{Insert}(J(a, b), s) \\ \mathbf{else} s) y) x \end{aligned}$$

For example,

$$\begin{aligned} \mathbf{join}(\text{employees}, \text{departments}, \\ \lambda(\text{emp}, \text{dept}).(\text{emp.dno}=\text{dept.dno} \mathbf{and} \text{dept.name}=\text{"CSE"}), \\ \lambda(\text{emp}, \text{dept}).\text{emp}) \end{aligned}$$

A function f is **commutative-idempotent** if:

$$\forall m \forall n \forall s : f(m, f(n, s)) = f(n, f(m, s)) \quad (\mathbf{commutativity})$$

$$\forall n \forall s : f(n, f(n, s)) = f(n, s) \quad (\mathbf{idempotence})$$

If a function f_s is commutative-idempotent then

$$\text{fold}^{set}(f_e, f_s)(\text{Insert}(a, s)) = f_s(a, \text{fold}^{set}(f_e, f_s) s)$$

Promotion Theorem for Sets:

$$\phi_e() = g(f_e())$$

$$\phi_s(a, g(r)) = g(f_s(a, r))$$

$$g(\text{fold}^{set}(f_e, f_s) x) = \text{fold}^{set}(\phi_e, \phi_s) x$$

and if f_s is commutative-idempotent then so is ϕ_s .

The Term Language

A term τ in the language has one of the following forms:

- **variable:** x ;
- **construction:** $C(\tau_1, \dots, \tau_n)$;
- **fold:** $\text{fold}^T(f_1, \dots, f_n) \tau$,
where each f_i has the form $\lambda(x_1, \dots, x_m). \tau_i$.

Safe Term: A program in the term language is *safe* if it does not contain terms $\text{fold}^T(\bar{f}) t$, where t contains a reference to an *accumulative result variable*.

e.g. this is *safe*:

$$\text{append}(x, y) = \text{fold}^{list}(\lambda().y, \lambda(a, r).\text{Cons}(a, r)) x$$

e.g. this is *NOT safe*:

$$\begin{aligned} \text{reverse}(x) &= \text{fold}^{list}(\lambda().\text{Nil}, \lambda(a, r).\text{append}(r, \text{Cons}(a, \text{Nil}))) x \\ &= \text{fold}^{list}(\lambda().\text{Nil}, \\ &\quad \lambda(a, r).\text{fold}^{list}(\lambda().\text{Cons}(a, \text{Nil}), \\ &\quad \quad \lambda(b, s).\text{Cons}(b, s)) r) x \end{aligned}$$

Canonical Terms

A **canonical term** is a term in which

- all folds in the term are over variables;
- none of these variables are accumulative result variables.

Theorem: Any safe term can be transformed into a canonical form.

The Normalization Algorithm (Example)

We will improve $\text{length}(\text{append}(x, y))$, where:

$$\text{length}(x) = \text{fold}^{list}(\lambda().\text{Zero}, \lambda(a, r).\text{Succ}(r)) x$$

$$\text{append}(x, y) = \text{fold}^{list}(f_n, f_c) x \quad \text{where} \quad \begin{cases} f_n = \lambda().y \\ f_c = \lambda(a, r).\text{Cons}(a, r) \end{cases}$$

We need to find some

$$\text{fold}^{list}(\phi_n, \phi_c) x = \text{length}(\text{append}(x, y))$$

The *list* promotion theorem is:

$$\begin{array}{l} 1) \phi_n() = g(f_n()) \\ 2) \phi_c(a, g(r)) = g(f_c(a, r)) \end{array}$$

$$g(\text{fold}^{list}(f_n, f_c) x) = \text{fold}^{list}(\phi_n, \phi_c) x$$

We apply this theorem with $g = \text{length}$:

$$\begin{aligned} 1) \phi_n() &= \text{length}(f_n()) \\ &= \text{length}(y) \end{aligned}$$

$$\begin{aligned} 2) \phi_c(a, \text{length}(r)) &= \text{length}(f_c(a, r)) \\ &= \text{length}(\text{Cons}(a, r)) \\ &= \text{Succ}(\text{length}(r)) \end{aligned}$$

$$\Rightarrow \phi_c(a, u) = \text{Succ}(u)$$

Therefore:

$$\text{length}(\text{append}(x, y)) = \text{fold}^{list}(\lambda().\text{length}(y), \lambda(a, u).\text{Succ}(u)) x$$

The Normalization Algorithm

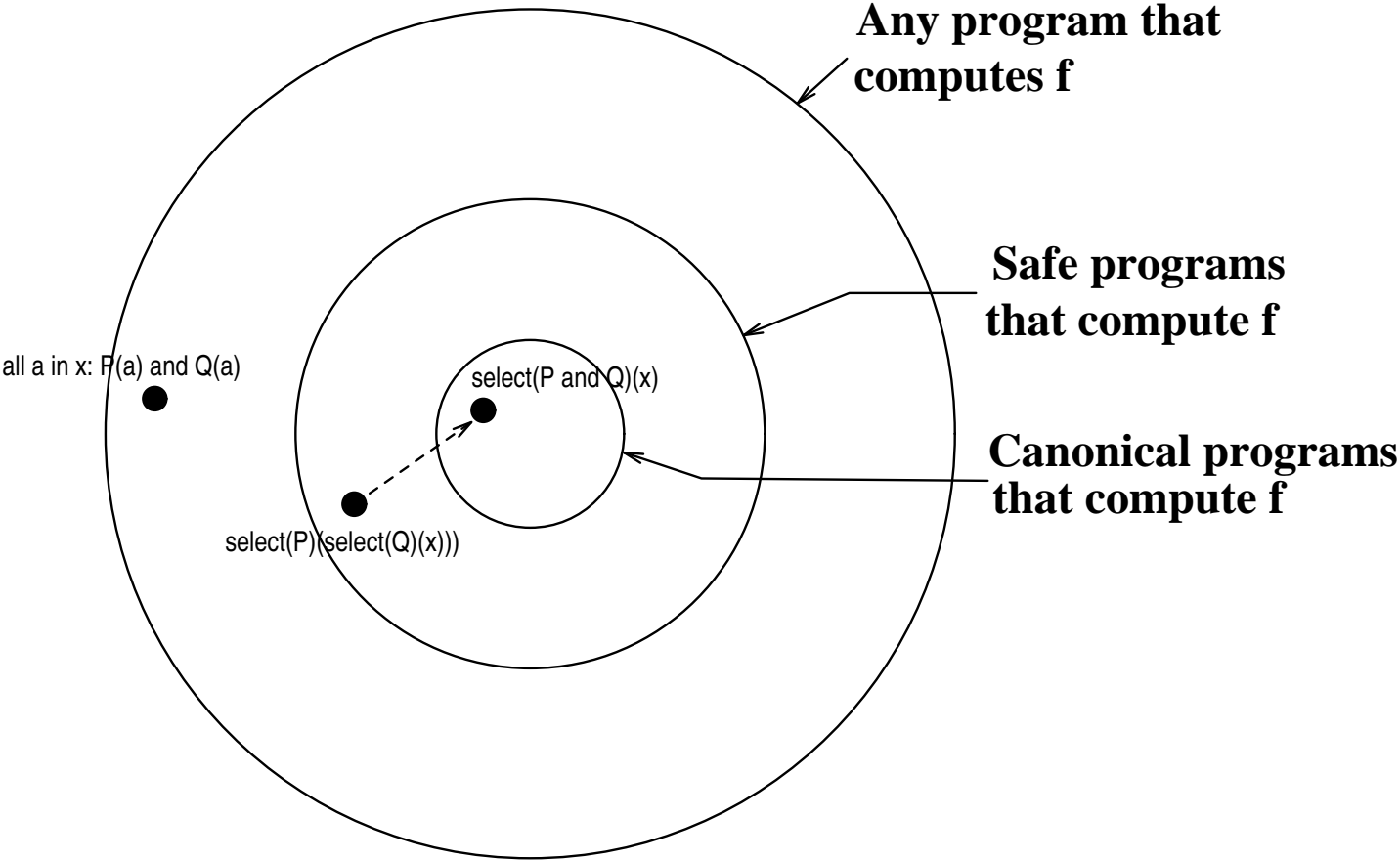
The normalization algorithm is a calculation-based, $\mathcal{O}(n \log n)$, algorithm that fuses piped folds:

$$\text{fold}^S(\bar{g}) (\text{fold}^T(\bar{f}) x) \longrightarrow \text{fold}^T(\bar{\phi}) x$$

It facilitates the following:

- automation of the unfold-simplify-fold method;
- automation of the techniques of deforestation and loop fusion;
- generalization of many well known algebraic optimizations
e.g. $\text{map}(f) \circ \text{map}(g) = \text{map}(f \circ g)$
and $\text{restrict}(p_1) \circ \text{restrict}(p_2) = \text{restrict}(p_1 \wedge p_2)$;
- automation of a form of partial evaluation;
- inductive theorem proving;
- implementation of most stream-based pipelined techniques during optimization time (instead of during evaluation time).

The Fold Optimization Algorithm



The **fold optimization algorithm** is a search-based algorithm that enumerates all canonical programs that satisfy an equation:

e.g. find f in $g \circ f = \phi$ given g and ϕ .

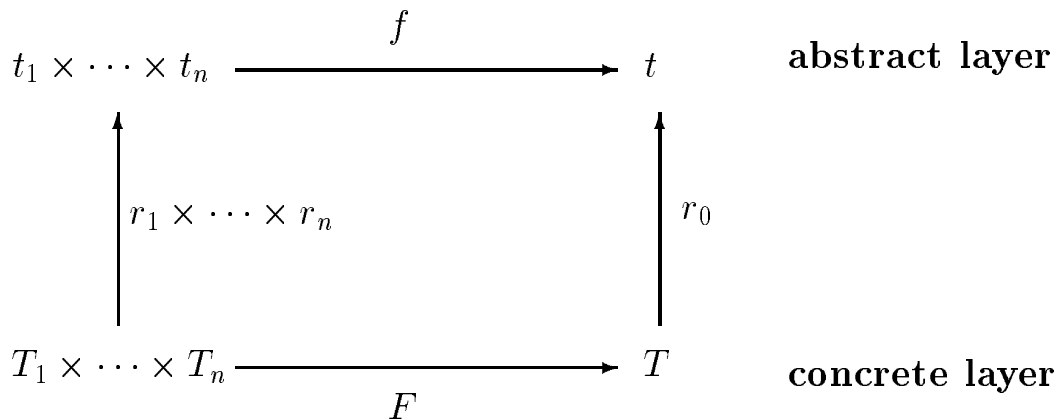
- it performs program synthesis based on second-order pattern matching;
- most alternatives are derived from commutative operations (not from associative);
- the search tree is a long (sometimes infinite) but not a bushy tree (max branch factor is 3);
- it can be guided by heuristics and cost functions.

Type Transformation

Database implementation is the translation of *abstract functions* or queries that operate on abstract values into efficient *concrete algorithms* that manipulate storage structures.

The type transformation model:

- the designer specifies how abstract types are mapped into storage structures by providing **abstraction functions**;
- the translator uses this information to translate abstract functions into concrete programs.



The concrete function F is an *implementation* of the abstract function f iff:

$$r_0 \circ F = f \circ (r_1 \times \cdots \times r_n)$$

Type transformation supports a high degree of *data independence*: abstract values and their operations can be significantly independent of their implementations.

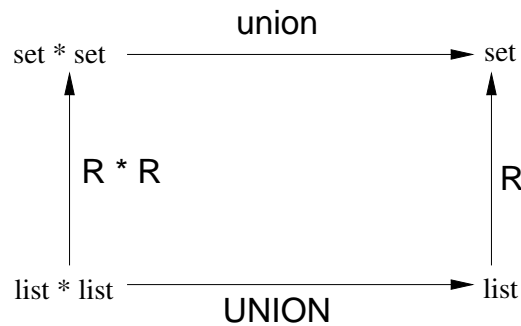
Example: Translating Set Operations

Abstract Operation:

$$\text{union}(x, y) = \text{fold}^{\text{set}}(\lambda().y, \lambda(a, r).\text{Insert}(a, r)) x$$

Abstraction function:

$$R = \text{fold}^{\text{list}}(\lambda().\text{Emptyset}, \lambda(a, r).\text{Insert}(a, r))$$



Any implementation UNION of union must satisfy:

$$R \circ \text{UNION} = \text{union} \circ (R \times R)$$

or equivalently:

$$\forall x, y : R(\text{UNION}(x, y)) = \text{union}(R(x), R(y))$$

e.g. one solution for UNION is list append.

Conclusion

- folds are uniformly defined over a wide spectrum of bulk data structures, which includes sets, lists, trees, tuples, numbers, and booleans;
- canonical programs are expressive enough to capture a large class of queries;
- the explicit iteration structure of fold programs supports calculation-based algebraic optimizations;
- there are very few canonical programs that compute a function, so search for the best is feasible in most cases;
- the fold optimization algorithm is a search-based algorithm that explores the reduced space of equivalent canonical forms.