

Deriving Incremental Programs

Yanhong Liu*
Tim Teitelbaum*

TR 93-1384
September 1993

Department of Computer Science
Cornell University
Ithaca, NY 14853-7501

*The authors gratefully acknowledge the support of the Office of Naval Research under contract No. N00014-92-J-1973.

Deriving Incremental Programs

Yanhong Liu* Tim Teitelbaum*
Department of Computer Science, Cornell University, Ithaca, NY 14853
{yanhong, tt}@cs.cornell.edu

September 1993

Abstract

A systematic transformational approach is given for deriving incremental programs from non-incremental programs. We exploit partial evaluation, other static analysis and transformation techniques, and domain-specific knowledge in order to provide a degree of incrementality not otherwise achievable by a generic incremental evaluator. Illustrative examples using the transformation approach are given.

1 Introduction

Incremental programs take advantage of repeated computations on inputs that differ only slightly from one another, avoiding unnecessary duplication of common computations. Given a computation of $f(x)$ with result r and a new input $x \oplus \delta x$ to the program f , an incremental program f' computes the result of $f(x \oplus \delta x)$ by making use of r , avoiding computations previously done by $f(x)$. Methods of incremental computation have widespread application throughout software, e.g., optimizing compilers [1] [6] [7], transformational programming [15] [32], interactive editing systems [25] [3], etc. Based on the observation that explicitly incremental programs are hard to write, we aim to provide incremental computation automatically from non-incremental programs.

In this paper, we present a technique for formally deriving an incremental program f' from a non-incremental program f and input change \oplus . We exploit partial evaluation, other static analysis and transformation techniques, and domain-specific knowledge in an attempt to provide a degree of incrementality not otherwise achievable by a generic incremental evaluator. The technique is presented as a four-step transformation: (1) expanding the computation $f(x \oplus \delta x)$ to separate computations on x from the rest, (2) introducing the cache parameter $r = f(x)$, (3) replacing

*The authors gratefully acknowledge the support of the Office of Naval Research under contract No. N00014-92-J-1973

equivalent computations on x by retrieving results from r , and (4) eliminating redundant code. In the third step, a subroutine called auxiliary generalized partial evaluation is employed to help discover incrementalities. We also cast our derivation approach in terms of a generalized partial evaluation point of view.

Paige’s work on finite differencing is the pioneering work on deriving incremental programs [18] [14] [31] [5]. Whereas Paige’s methods are for very high-level set-theoretic languages, our techniques apply to general functional programs, as does Smith’s approach in KIDS [32] [33]. For further discussions and comparisons between these methods, see Section 8.

The rest of the paper is organized as follows: In Section 2, we give the definition of incremental programs. Our basic approach for the derivation of incremental programs is outlined in Section 3. The formalization of the problem with the language we are using is presented in Section 4. In Section 5, we formalize the transformations for the derivation and argue that they preserve the semantics of the programs. An example using the transformations is given in Section 6, in which an insert program is derived from a sort program. Several related issues, including a limitation of the approach and directions for remedying them, are discussed in Section 7. Finally, we discuss related work in incremental computation in Section 8.

2 Incremental Programs

Let f be a program and x an input for f . Let δx be a *change* in the value of x and \oplus a *change operation* that combines value x and change δx to produce a new input value $x \oplus \delta x$. In general, the domains of x and δx need not be the same. For example, consider a list of integers x , an integer δx , and $x \oplus \delta x = \text{cons}(\delta x, x)$; i.e., x is changed by prepending integer δx at the head of the list.

Let $r = f(x)$; i.e., r is the result of executing program f on input x . We call r the *cached result* and aim to make use of it in computing $f(x \oplus \delta x)$. Let f'_{\oplus} be a program such that whenever r is the cached result $f(x)$, $f'_{\oplus}(x, \delta x, r)$ computes $f(x \oplus \delta x)$ for all possible x and δx , i.e.,

$$(\forall x)(\forall \delta x)(\forall r) [f(x) \downarrow = r \implies f'_{\oplus}(x, \delta x, r) = f(x \oplus \delta x)], \quad (1)$$

where $f(x) \downarrow = r$ means that $f(x)$ terminates with value r . Then f'_{\oplus} is called an *incremental version* of f under \oplus if, in addition to satisfying (1), $f'_{\oplus}(x, \delta x, r)$ computes faster than $f(x \oplus \delta x)$ for almost all x and δx and makes non-trivial use of r . By computing faster, we mean that a computation takes less time under some time model. By making non-trivial use of r , we mean that the availability of r is indispensable in obtaining the faster computation.

For example, assume the RAM model and let $\text{sort}(x)$ be a selection sort on a list x of numbers

that uses a naive select (which finds a least element in a list by traversing the list from the beginning to the end). Then, $sort(x)$ takes time $O(n^2)$ for x of length n . Consider the change to the input x given by $x \oplus \delta x = cons(\delta x, x)$, where δx is a number. There exists a $sort'(x, \delta x, r)$ that, when r is $sort(x)$, computes $sort(cons(\delta x, x))$ by inserting δx in r in $O(n)$ time. Then by definition, $sort'$ is an incremental version of $sort$. In particular, $sort'$ makes non-trivial use of r since otherwise it is impossible for a sort program to compute in $O(n)$ time. In contrast, consider a $sort''(x, \delta x, r)$ that computes $sort(cons(\delta x, x))$ by doing a merge sort on $cons(\delta x, x)$ and takes only $O(n \log n)$ time. $sort''$ is not an incremental version of $sort$ since it gains its speedup without the use of r .

Suppose f'_\oplus is an incremental version of f under \oplus . We say that f'_\oplus is a *complete incremental version* if f'_\oplus does not depend at all on the parameter x , i.e., $f'_\oplus(x, \delta x, r) = f'_\oplus(c, \delta x, r)$ for all $x, \delta x$, and constant c . In other words, $f'_\oplus(x, \delta x, r)$ computes $f(x \oplus \delta x)$ but only uses the change parameter δx and the cached result r . In the sort example, $sort'$ is a complete incremental version of $sort$ under \oplus .

For a formal definition of incremental programs, see [13]. When no confusion arises, we simply write f' for f'_\oplus . For typographical convenience, we shall typically use y for the change parameter rather than δx . Also, throughout the paper, identifier r always refers to the cached result of a previous computation $f(x)$.

In [13], it is proved that in general, given a program f and an input change operation \oplus , whether an incremental version f' exists is undecidable. Accordingly, we compromise by seeking an approximation to an incremental version, which we call a *differentiated version*. A differentiated version f' differs from an incremental version in that it only guarantees that f' computes at least as fast as the original program for almost all inputs, while an incremental version computes strictly faster than the original program for almost all inputs. If a differentiated version only uses the change parameter δx and the cached result r , then it is called a *complete differentiated version*.

Obviously, a differentiated version $f'(x, \delta x, r)$ can be trivially defined to be $f(x \oplus \delta x)$. But then no efficiency is gained by f' , which is of no interest. The goal is to make f' as efficient as possible by having it use the cached result r as much as possible.

3 Derivation Approach

Basic Ideas. We aim to identify in the computation of $f(x \oplus y)$ those computations that are also performed in the computation of $f(x)$ and avoid recomputing them if possible.

More precisely, by expanding the computation of $f(x \oplus y)$, those computations on x in $f(x \oplus y)$

can be separated out from the rest of the computation. Given the cached result r of a previous computation of $f(x)$, the results of some of the computations on x in $f(x \oplus y)$ may be retrieved more efficiently from r than computing them from scratch. This more efficient way to compute $f(x \oplus y)$ is captured in the definition of the differentiated version $f^*(x, y, r)$.

In order to explain the basic ideas further, we break the derivation approach into four steps.

Step 1. Expanding the computation of $f(x \oplus y)$ to separate out the computations on x from the rest of computation. This is achieved by unfolding and simplifying function applications whose arguments involve both x and y .

Step 2. Introducing the cached result r of the computation $f(x)$. In this step, we define $f^*(x, y, r)$ to be the expanded computation of $f(x \oplus y)$ obtained from Step 1. In particular, if an application of f whose arguments involve both x and y occurs recursively in the expanded computation of $f(x \oplus y)$, it is replaced by an instance of $f^*(x, y, r)$, if possible. Thus, the computation of $f(x \oplus y)$ is performed by $f^*(x, y, r)$.

Step 3. Using the cached result r in the computation of $f^*(x, y, r)$. In the definition of $f^*(x, y, r)$ obtained from Step 2, we find computations that are also performed by $f(x)$, and replace them by retrievals from the cached result r where possible. As a straightforward example, an occurrence of $f(x)$ in the definition of $f^*(x, y, r)$ can be replaced by r . Moreover, within the context of a subexpression in the definition of $f^*(x, y, r)$, the computation $f(x)$ may have a specialized form. If the subexpression has that specialized form, then it can also be replaced by r .

Step 4. Eliminating redundant computations in the resulting definition obtained from Step 3. This can be done using standard optimizations. In particular, since some computations depending on x have been replaced by computations on r , part or all of parameter x may have become redundant. Eliminating them results in a more efficient differentiated version $f'(\bar{x}, y, r)$, where \bar{x} is x or the part of x that remains as a parameter.

Example. To illustrate the four-step derivation approach, we present a simple case of matrix multiplication as an example. Let x be $\langle C, R \rangle$, where C and R are two lists of numbers. The function $mtxMul(C, R)$ treats C as a column matrix, R as a row matrix, and returns the product of the two matrices. The resulting matrix is represented as a list of rows, where each row is a list of numbers. The auxiliary function $rowMul(e, R)$ multiplies an element e with each element in a

row R and returns the product row. The function definitions are:

$$\begin{aligned}
mtxMul(C, R) &= \text{if } null(C) \text{ then } nil \\
&\quad \text{else } cons(rowMul(car(C), R), mtxMul(cdr(C), R)) \\
rowMul(e, R) &= \text{if } null(R) \text{ then } nil \\
&\quad \text{else } cons(e * car(R), rowMul(e, cdr(R)))
\end{aligned} \tag{2}$$

Suppose r is the cached result $mtxMul(C, R)$ and the new input to $mtxMul$, $x \oplus y$, is defined as $\langle C, R \rangle \oplus y = \langle C, cons(y, R) \rangle$. We aim to compute $mtxMul(C, cons(y, R))$ more efficiently by making use of r .

Step 1 considers

$$mtxMul(C, cons(y, R)).$$

First, we unfold $mtxMul(C, cons(y, R))$ since its arguments involve both x and y . We get

$$\begin{aligned}
&\text{if } null(C) \text{ then } nil \\
&\text{else } cons(rowMul(car(C), cons(y, R)), mtxMul(cdr(C), cons(y, R))).
\end{aligned} \tag{3}$$

Then, we unfold $rowMul(car(C), cons(y, R))$ in (3) since its arguments also involve both x and y . We get

$$\begin{aligned}
&\text{if } null(C) \text{ then } nil \\
&\text{else } cons(\text{if } null(cons(y, R)) \text{ then } nil \\
&\quad \text{else } cons(car(C) * car(cons(y, R)), rowMul(car(C), cdr(cons(y, R)))), \\
&\quad mtxMul(cdr(C), cons(y, R))).
\end{aligned} \tag{4}$$

Using the axioms of data types $null(cons(h, t)) = false$, $car(cons(h, t)) = h$, and $cdr(cons(h, t)) = t$ and specializing the false conditional expression in (4) to its false branch, we get

$$\begin{aligned}
&\text{if } null(C) \text{ then } nil \\
&\text{else } cons(cons(car(C) * y, rowMul(car(C), R)), mtxMul(cdr(C), cons(y, R))).
\end{aligned} \tag{5}$$

The function application $mtxMul(cdr(C), cons(y, R))$ in (5) is not further unfolded since it is an instance of $mtxMul(C, cons(y, R))$, which has already been unfolded.

Step 2 defines $mtxMul^*(\langle C, R \rangle, y, r)$ based on (5). The expression $mtxMul(cdr(C), cons(y, R))$ in (5) is equivalent to $mtxMul(\langle cdr(C), R \rangle \oplus y)$. Replacing this application of $mtxMul$ with an application of $mtxMul^*$, we get

$$\begin{aligned}
&\text{if } null(C) \text{ then } nil \\
&\text{else } cons(cons(car(C) * y, rowMul(car(C), R)), \\
&\quad mtxMul^*(\langle cdr(C), R \rangle, y, mtxMul(cdr(C), R))),
\end{aligned} \tag{6}$$

which is then used as the initial definition of $mtxMul^*(\langle C, R \rangle, y, r)$. Note the invariance that the third parameter of $mtxMul^*$ is the application of $mtxMul$ to the first parameter of $mtxMul^*$.

Step 3 aims to use the cached result r to compute (6) more efficiently. In particular, in the false branch of the conditional expression, $null(C)$ is false. In this context, we specialize (2) and find that

$$mtxMul(C, R) = cons(rowMul(car(C), R), mtxMul(cdr(C), R)). \quad (7)$$

From the invariance $mtxMul(C, R) = r$, we have $cons(rowMul(car(C), R), mtxMul(cdr(C), R)) = r$, which implies,

$$rowMul(car(C), R) = car(r) \text{ and } mtxMul(cdr(C), R) = cdr(r). \quad (8)$$

After substitutions in (6) using the last two equations, we get

$$\begin{aligned} & \text{if } null(C) \text{ then } nil \\ & \text{else } cons(cons(car(C) * y, car(r)), mtxMul^*(\langle cdr(C), R \rangle, y, cdr(r))), \end{aligned} \quad (9)$$

which is a new definition of $mtxMul^*(\langle C, R \rangle, y, r)$.

Step 4 eliminates redundant code in the definition (9) of $mtxMul^*$. In particular, the parameter R of $mtxMul^*(\langle C, R \rangle, y, r)$ is not used as data or for any control, except in the recursive call to $mtxMul^*$ itself as the corresponding argument. After eliminating it, we get a final definition

$$mtxMul'(C, y, r) = \begin{aligned} & \text{if } null(C) \text{ then } nil \\ & \text{else } cons(cons(car(C) * y, car(r)), mtxMul'(cdr(C), y, cdr(r))). \end{aligned} \quad (10)$$

Therefore, given that $mtxMul(C, R) = r$, we have that $mtxMul'(C, y, r) = mtxMul^*(\langle C, R \rangle, y, r) = mtxMul(C, cons(y, R))$. The function $mtxMul'$ is a differentiated version of $mtxMul$ for the given notion of input change. Note that it is not a complete differentiated version since C can not be eliminated from (10).

Generalized Partial Evaluation. It is convenient and fruitful to view our derivation approach as an instance of *generalized partial evaluation* (GPE). GPE is a semantic-based program transformation: given an expression e and an information set I that describes partial knowledge about e , GPE produces a *residual expression* e' by specializing e using information from I .

The notion of GPE has appeared in [10] [37], where it was called generalized partial computation (GPC). It was originally proposed to overcome the drawbacks of conventional partial evaluation, which only specializes programs given static values for certain arguments. The key idea of GPE

is to employ a theorem prover during the specialization to utilize information such as the logical structure of programs, algebraic properties of primitive functions, and axioms for abstract data types, etc.

Expanding the computation of $f(x \oplus y)$ in Step 1 can be regarded as a special case of GPE in which the expression e is $f(z)$ and the information set states that z is $x \oplus y$ and the cached result $f(x)$ is available. We call this use of GPE the *original generalized partial evaluation* (oGPE); the residual expression of the oGPE is called the *original residual expression*. In the matrix multiplication example above, (5) is the original residual expression.

The specialization of $f(x)$ in a specific context in Step 3 can also be regarded as a special case of GPE. We call this the *auxiliary generalized partial evaluation* (aGPE); a residual expression of the aGPE is called an *auxiliary residual expression*. For the example above, the specialization of $mtxMul(C, R)$ in the context where $null(C)$ is false is an aGPE that yields the right hand side of (7) as an auxiliary residual expression.

In [10] and [37], GPE transformation methods for a restricted form of language (u-form) and a lazy functional language are discussed, and small examples are given. In this paper, we see that combining the idea of GPE transformations with other special transformation techniques results in an effective transformation approach for deriving incremental programs.

4 Formalization of the Problem

First, we present the syntax of the simple language we use and give it a simple evaluation semantics. Then, the derivation problem is described in terms of the language.

4.1 Language

We consider a pure first-order functional language with the following syntax.

Syntactic Domains

| | |
|-----------|---|
| Var | variables, ranged over by x, x_0, \dots |
| Con | constructors, ranged over by c, c_0, \dots |
| $PrimFun$ | primitive functions, ranged over by p, p_0, \dots |
| Fun | functions, ranged over by f, f_0, \dots |
| Exp | expressions, ranged over by e, e_0, \dots |

Abstract Syntax

| | | |
|----------------------|-------------|-------------------------|
| For $e \in Exp$, | | |
| $e ::= v$ | $v \in Var$ | variable |
| $c(e_1, \dots, e_n)$ | $c \in Con$ | constructor application |

| | | | |
|--|---|-----------------|--------------------------------|
| | $p(e_1, \dots, e_n)$ | $p \in PrimFun$ | primitive function application |
| | $f(e_1, \dots, e_n)$ | $f \in Fun$ | function application |
| | if e_1 then e_2 else e_3 | | conditional expression |
| | let $v = e_1$ in e_2 end | | binding expression |

A program is a set F of mutually recursive function definitions

$$\{f_i(x_1, \dots, x_{n_i}) = e_i \mid i \in I_F\} \quad (11)$$

where I_F is a finite index set satisfying $0 \in I_F$. The function f_0 with parameter x_1, \dots, x_{n_0} is the expression to be evaluated in the context of these definitions. Note that, from this section on, f_0 refers to the function for which we aim to derive a differentiated version, whereas f refers to an arbitrary function in the set F . We assume that each v in a **let** expression is renamed so that it is distinct from all other variables used in F .

A simple strict semantics for expressions can be given by semantic function $\mathcal{E} : Exp \rightarrow Env \rightarrow Val$, where Env and Val are semantic domains. In particular, constructor applications and primitive function applications are strict. Val is the set of values, including constants and constructions. Integers and Boolean values (T for true and F for false) are included among the constants. Env is a type-respecting mapping from variables to values. Formally, we give the semantic domains and the semantic equations below, where D_\perp is the domain D extended with a value \perp that is below all elements in D and $v \mapsto d$ means that the variable v is bound to the value d . Sometimes, we use $[x \mapsto d]$ as an abbreviation for $[x_1 \mapsto d_1, \dots, x_n \mapsto d_n]$.

Semantic Domains

$Val = Constants + Constructions$, ranged over by d

where $Constants = (Con)_\perp$,

$Constructions = (Con \times Val)_\perp + (Con \times Val \times Val)_\perp + \dots + (Con \times Val^n)_\perp + \dots$

$Env = Var \rightarrow Val$, ranged over by ρ

Semantic Equations

$$\begin{aligned}
\mathcal{E}[[v]]\rho &= \rho(v) \\
\mathcal{E}[[c(e_1, \dots, e_n)]]\rho &= c(\mathcal{E}[[e_1]]\rho, \dots, \mathcal{E}[[e_n]]\rho) \\
\mathcal{E}[[p(e_1, \dots, e_n)]]\rho &= p(\mathcal{E}[[e_1]]\rho, \dots, \mathcal{E}[[e_n]]\rho) \\
\mathcal{E}[[f(e_1, \dots, e_n)]]\rho &= \mathcal{E}[[e]]\rho[x_1 \mapsto d_1, \dots, x_n \mapsto d_n], \text{ if } d_i \text{ is of the right type for } x_i \text{ of } f \\
&= \perp, \text{ otherwise,} \\
&\text{where } d_i = \mathcal{E}[[e_i]]\rho, \text{ for } i = 1..n, \text{ and } f \text{ is defined by } f(x_1, \dots, x_n) = e \\
\mathcal{E}[[\text{if } e_1 \text{ then } e_2 \text{ else } e_3]]\rho &= \mathcal{E}[[e_2]]\rho, \text{ if } d = T \\
&= \mathcal{E}[[e_3]]\rho, \text{ if } d = F \\
&= \perp, \text{ otherwise,} \\
&\text{where } d = \mathcal{E}[[e_1]]\rho \\
\mathcal{E}[[\text{let } v = e_1 \text{ in } e_2 \text{ end}]]\rho &= \mathcal{E}[[e_2]]\rho[v \mapsto d], \text{ if } d \text{ of the right type for } v, \\
&= \perp, \text{ otherwise,} \\
&\text{where } d = \mathcal{E}[[e_1]]\rho
\end{aligned}$$

4.2 Description of the Problem

We are given a set F of function definitions, as in (11), and an *input change operation* \oplus . The operation \oplus describes the change to the arguments of the function f_0 . It combines any tuple of old arguments $x \equiv \langle x_1, \dots, x_n \rangle^{1,2}$ and a tuple of variables $y \equiv \langle y_1, \dots, y_m \rangle$ to form a tuple of new arguments $x' \equiv \langle x'_1, \dots, x'_n \rangle$, i.e., $x' = x \oplus y$. The operation \oplus can be regarded as a tuple of n separate *single change operations* $\oplus_1, \dots, \oplus_n$ for x'_1, \dots, x'_n respectively, such that $x'_i = x \oplus_i y$ for each $i = 1, \dots, n$. Each of the single change operations \oplus_i is defined in such a way that x'_i is computed as some function g_i of parameters x_j 's and y_k 's. Usually, g_i is a constructor application or a primitive function application. The most typical y and \oplus are the tuple $y \equiv \langle y_1, \dots, y_n \rangle$ and the operation \oplus such that there exists a binary function g_i of x_i and y_i , such that $x'_i = g_i(x_i, y_i)$, for each $i = 1, \dots, n$.

Note that in order to keep our language simple, we do not include tuples as expressions. However, tuples are used in the presentation as a notational convenience. In particular, if $x \equiv \langle x_1, \dots, x_n \rangle$, $y \equiv \langle y_1, \dots, y_m \rangle$, $\oplus \equiv \langle \oplus_1, \dots, \oplus_n \rangle$, and $x \oplus_i y$ is defined using a function g_i of x_j 's and y_k 's, then $f_0(x)$ is an abbreviation for $f_0(x_1, \dots, x_n)$, $f_0^*(x, y, f_0(x))$ is an abbreviation for $f_0^*(x_1, \dots, x_n, y_1, \dots, y_m, f_0(x_1, \dots, x_n))$, and $x \oplus y$ is an abbreviation for $\langle x \oplus_1 y, \dots, x \oplus_n y \rangle$, where $x \oplus_i y$ is an abbreviation for $g_i(x_1, \dots, x_n, y_1, \dots, y_m)$. These abbreviations are used not only for tuples of variables, but also for tuples of expressions.

We are only interested in using the cached result if we can save time by doing so. Accordingly, we must have a model that describes the time needed to compute expressions. We use $T(e)$ to denote the time needed to compute expression e . If r is the only free variable in an expression e , then $T(e)$ is described as “the time needed to retrieve the value of e from cached result r ”. If $e \equiv \langle e_1, \dots, e_n \rangle$, we use $T(e)$ to denote the sum of the $T(e_i)$'s. The function T can be obtained using fairly standard constructions [38] [27] [30]. Note that different primitive functions may take different time to evaluate.

In general, given two expressions e_1 and e_2 that compute the same function of x , it is not decidable whether e_1 computes faster than e_2 for a given value of x . Therefore, we say that $T(e_1) \leq T(e_2)$ (or $T(e_1) < T(e_2)$) for a given value of x if we can *effectively* confirm the inequality using heuristics for the function T . We write $t(e_1) \leq t(e_2)$ to denote that we can effectively decide that there exists a constant c such that for all x , $T(e_1) \leq c \cdot T(e_2)$. In particular, we use $t(e_1) \leq t(e_2)$ to denote that for all x , $T(e_1) \leq T(e_2)$. Notice that $t(e_1) \leq t(e_2)$ (or $t(e_1) \leq t(e_2)$) is just notation

¹In the following presentation, we use ‘ \equiv ’ for identity.

²For simplicity, the arity n_0 of function f_0 referred to in (11) is denoted by n .

as t itself is not defined.

As a matter of fact, we do not need to know the exact time needed to evaluate an expression. Instead, we only care about comparing the times needed to evaluate two expressions. Therefore, we can use heuristics to compare times. For example, if $T(e) > T(e_1) + \max\{T(e_2), T(e_3)\}$, then $T(e) > T(\text{if } e_1 \text{ then } e_2 \text{ else } e_3)$. If $T(e_1) \leq T(e_2)$, then $T(p(e_1)) \leq T(p(e_2))$ for any applicable primitive function p . In particular, if e_1 is a subexpression of e_2 , then $T(e_1) \leq T(e_2)$. For example, $T(e_1) \leq T(\text{car}(\text{cons}(e_1, e_2)))$. We can also make use of the transitivities of various comparison relations.

Given F , \oplus , and the time model T , we aim to apply a series of transformations to obtain a finite set F' of function definitions $\{f'_i(x'_1, \dots, x'_{n'_i}) = e'_i \mid i \in I_{F'}\}$, where $0 \in I_{F'}$ and f'_0 is a differentiated version of f_0 under \oplus .

5 Formalization of the Transformations for the Derivation

In order to make the formalization of the transformation approach tractable, we break it into four steps:

| step | description | transformation |
|------|---|------------------------|
| 1 | expanding $f_0(x \oplus y)$ to separate computations on x and y | \mathcal{G}_o |
| 2 | introducing the cache parameter r | \mathcal{D}_t |
| 3 | using the cache parameter | \mathcal{S}_t |
| 4 | eliminating redundant code | standard optimizations |

The definition of a differentiated version f_0^* of f_0 is obtained by transforming $f_0(x \oplus y)$ with \mathcal{G}_o , \mathcal{D}_t , and \mathcal{S}_t in turn. The final differentiated version f'_0 is obtained by applying standard optimizations, especially redundant code elimination, on f_0^* and the other function definitions in the set F . In the following, we discuss each step separately and argue (a) that together they preserve the semantics of the program, and (b) that computations using the derived programs are at least as fast as using the original programs. An example using the transformations is given in Section 6.

Since this formalization breaks the transformation into separate steps, it has some obvious limitations and inefficiencies. For a characterization of this simple formalization and ideas for a more general approach, see Section 7.

5.1 Information Sets

Section 3 introduced the GPE view of our derivation approach. GPE specializes an expression e with respect to an information set I . We use $Info$ to denote the domain of information sets. For

any $I \in \text{Info}$, I is a set of equations whose conjunction represents the information with which an expression is to be specialized. An equation is a pair of expressions and has type $\text{Exp} \times \text{Exp}$, where the two expressions are related by the Boolean primitive equality function. In this paper, we use $e \leftrightarrow e'$ to denote that e equals e' for two expressions e and e' . Usually, an information set is collected from the context of an expression. For example, let e be the expression

$$\text{let } v = p(x) \text{ in if } v < x \text{ then } e_1 \text{ else } e_2 \text{ end.}$$

The information set for subexpression e_2 in e is $\{v \leftrightarrow p(x), v < x \leftrightarrow F\}$.

5.2 Underlying Logic

The key idea of GPE is to use a theorem prover to make inferences based on its logic and the information collected about a program. We are only interested in logics that contain axioms, inference rules, and basic formulas that are true for the language we are using. For any information set I , let e_I be a Boolean-valued expression in our language that represents the information of the set I . A logic is called an *underlying logic* [10] for our language if it is *compatible* with the evaluation semantics of our language, where compatibility is defined as follows:

Definition 5.1 *A logic L is compatible with the evaluation semantics \mathcal{E} of our language if and only if, for any two information sets I and I' such that the equations in I' are provable from those in I based on the logic L , we have $\forall v_z, \mathcal{E}[[e_I]] [z \mapsto v_z] = T \implies \mathcal{E}[[e_{I'}]] [z \mapsto v_z] = T$, where z_1, \dots, z_k are all the free variables in e_I and $e_{I'}$.*

In the transformations below, we assume that L_0 is some fixed underlying logic of our language. In principle, L_0 can be any underlying logic, depending on the power we desire for making inferences about our programs. In this paper, we assume that a theorem prover based on the logic L_0 is available, and we write $e \leftrightarrow_I^* e'$ to denote that a finite proof that e equals e' can be found by the theorem prover using equations in the set I .

5.3 Step 1: Separating x and y

Step 1 specializes the function f_0 for arguments of the form $x \oplus y$. The goal of this step is to expand computations of $f_0(x \oplus y)$ into computations on x and y separately in the hope that computations on x alone can be avoided by making use of the cached result $f_0(x)$.

Step 1 unfolds function applications whose arguments involve computations on both x and y . Conversely, function applications whose arguments depend purely on x or y are not unfolded. The intuition behind this strategy is as follows: First, by unfolding and simplifying function applications

on arguments involving both x and y , we can hope to expose computations on x and on y separately. Secondly, computations purely related to y are new and need to be carried out anyway. The only purpose in unfolding such function applications would be to optimize them, but incrementalization, not optimization, is the subject of this paper. Thirdly, the intention in unfolding function applications that depend only on x would be either to enable optimizations of such computations or to help equate them with computations performed by $f_0(x)$. The former is not the subject of this paper, as just mentioned. The latter is actually a theorem proving problem, which is not the subject of this paper either. There is a fourth category of function applications, namely those that depend on neither x nor y . Such computations allow a straightforward optimization, i.e., evaluation by unfolding and simplification. However, this is a potential source of nontermination. We may wish to omit such optimizations, as will be discussed later for the unfolding rule. In short, our main concern is incrementality of f_0 with respect to the new input $x \oplus y$.

Step 1 is implemented by the transformation \mathcal{G}_o defined in Figure 1, where the type of \mathcal{G}_o is $Exp \rightarrow Info \rightarrow Exp$. Given an expression $e \in Exp$ and an information set $I \in Info$, \mathcal{G}_o specializes e with respect to I . By applying \mathcal{G}_o to the expression $f_0(x \oplus y)$ and the null information set \emptyset , we obtain expression e^o :

$$e^o \equiv \mathcal{G}_o[[f_0(x \oplus y)]] \emptyset. \quad (12)$$

The body of function \mathcal{G}_o consists of a set of cases conditioned on the structure of its argument expression. The first column in Figure 1 gives a name for each case, where a description of the case is encoded in the name:

| name | top-level construct |
|-------|--------------------------------|
| v | variable |
| c | constructor application |
| p | primitive function application |
| f | function application |
| if | if -expression |
| let | let -expression |

Subscript i in a name indicates that subexpression i within the named construct is being reduced. Subscript $i-if$ or $i-let$ indicates that the subexpression being reduced is an **if** or a **let** expression. Subscript e indicates that an expression itself is reduced after all of its subexpressions have been reduced. Subscript n indicates that an expression is not further reducible. The last column gives side conditions for the reduction rules.

The presentation of \mathcal{G}_o is simplified by omitting detailed control structures that sequence \mathcal{G}_o through its argument subexpressions. Thus, we just present the case of \mathcal{G}_o working on the

| Name | Transformation | Condition |
|---------------|--|--|
| v | $\mathcal{G}_o[[v]] I$ | |
| v_c | $= c$ | if $v \leftrightarrow_I^* c$ |
| v_n | $= v$ | otherwise |
| c | $\mathcal{G}_o[[c(e_1, \dots, e_n)]] I$ | |
| c_i | $= \mathcal{G}_o[[c(e_1, \dots, e_{i-1}, \mathcal{G}_o[[e_i]]I, e_{i+1}, \dots, e_n)]] I$ | if e_1, \dots, e_{i-1} are irreducible, not if or let , but e_i is reducible |
| c_{i-f} | $= \mathcal{G}_o[[\text{if } e'_1 \text{ then } c(e_1, \dots, e_{i-1}, e'_2, e_{i+1}, \dots, e_n) \text{ else } c(e_1, \dots, e_{i-1}, e'_3, e_{i+1}, \dots, e_n)]] I$ | if e_1, \dots, e_{i-1} are irreducible, not if or let , and e_i is irreducible and e_i is if e'_1 then e'_2 else e'_3 |
| c_{i-let} | $= \mathcal{G}_o[[\text{let } v = e'_1 \text{ in } c(e_1, \dots, e_{i-1}, e'_2, e_{i+1}, \dots, e_n) \text{ end}]] I$ | if e_1, \dots, e_{i-1} are irreducible, not if or let , and e_i is irreducible and e_i is let $v = e'_1$ in e'_2 end |
| c_n | $= c(e_1, \dots, e_n)$ | otherwise |
| p | $\mathcal{G}_o[[p(e_1, \dots, e_n)]] I$ | |
| p_i | $= \mathcal{G}_o[[p(e_1, \dots, e_{i-1}, \mathcal{G}_o[[e_i]]I, e_{i+1}, \dots, e_n)]] I$ | if e_1, \dots, e_{i-1} are irreducible, not if or let , but e_i is reducible |
| p_{i-f} | $= \mathcal{G}_o[[\text{if } e'_1 \text{ then } p(e_1, \dots, e_{i-1}, e'_2, e_{i+1}, \dots, e_n) \text{ else } p(e_1, \dots, e_{i-1}, e'_3, e_{i+1}, \dots, e_n)]] I$ | if e_1, \dots, e_{i-1} are irreducible, not if or let , and e_i is irreducible and e_i is if e'_1 then e'_2 else e'_3 |
| p_{i-let} | $= \mathcal{G}_o[[\text{let } v = e'_1 \text{ in } p(e_1, \dots, e_{i-1}, e'_2, e_{i+1}, \dots, e_n) \text{ end}]] I$ | if e_1, \dots, e_{i-1} are irreducible, not if or let , and e_i is irreducible and e_i is let $v = e'_1$ in e'_2 end |
| p_e | $= \mathcal{G}_o[[\text{Simp}[p(e_1, \dots, e_n)]]I] I$ | if e_1, \dots, e_n are irreducible, not if or let , but $p(e_1, \dots, e_n)$ can be simplified under I |
| p_n | $= p(e_1, \dots, e_n)$ | otherwise |
| f | $\mathcal{G}_o[[f(e_1, \dots, e_n)]] I$ | |
| f_i | $= \mathcal{G}_o[[f(e_1, \dots, e_{i-1}, \mathcal{G}_o[[e_i]]I, e_{i+1}, \dots, e_n)]] I$ | if e_1, \dots, e_{i-1} are irreducible, not if or let , but e_i is reducible |
| f_{i-f} | $= \mathcal{G}_o[[\text{if } e'_1 \text{ then } f(e_1, \dots, e_{i-1}, e'_2, e_{i+1}, \dots, e_n) \text{ else } f(e_1, \dots, e_{i-1}, e'_3, e_{i+1}, \dots, e_n)]] I$ | if e_1, \dots, e_{i-1} are irreducible, not if or let , and e_i is irreducible and e_i is if e'_1 then e'_2 else e'_3 |
| f_{i-let} | $= \mathcal{G}_o[[\text{let } v = e'_1 \text{ in } f(e_1, \dots, e_{i-1}, e'_2, e_{i+1}, \dots, e_n) \text{ end}]] I$ | if e_1, \dots, e_{i-1} are irreducible, not if or let , and e_i is irreducible and e_i is let $v = e'_1$ in e'_2 end |
| f_e | $= \mathcal{G}_o[[e[x_1 := e_1, \dots, x_n := e_n]] I$ where f is defined as $f(x_1, \dots, x_n) = e$ | if e_1, \dots, e_n are irreducible, not if or let , but the argument pattern is favorable and is not in the list l_f Also, put the argument pattern in the list l_f . |
| f_n | $= f(e_1, \dots, e_n)$ | otherwise |
| if | $\mathcal{G}_o[[\text{if } e_1 \text{ then } e_2 \text{ else } e_3]] I$ | |
| if_1 | $= \mathcal{G}_o[[\text{if } (\mathcal{G}_o[[e_1]]I) \text{ then } e_2 \text{ else } e_3]] I$ | if e_1 is reducible |
| if_{1-f} | $= \mathcal{G}_o[[\text{if } e'_1 \text{ then } (\text{if } e'_2 \text{ then } e_2 \text{ else } e_3) \text{ else } (\text{if } e'_3 \text{ then } e_2 \text{ else } e_3)]] I$ | if e_1 is irreducible and e_1 is if e'_1 then e'_2 else e'_3 |
| if_{1-let} | $= \mathcal{G}_o[[\text{let } v = e'_1 \text{ in } (\text{if } e'_2 \text{ then } e_2 \text{ else } e_3) \text{ end}]] I$ | if e_1 is irreducible and e_1 is let $v = e'_1$ in e'_2 end |
| if_2 | $= \mathcal{G}_o[[e_2]] I$ | if e_1 is irreducible and is not if or let , and $e_1 \leftrightarrow_I^* T$ |
| if_3 | $= \mathcal{G}_o[[e_3]] I$ | if e_1 is irreducible and is not if or let , and $e_1 \leftrightarrow_I^* F$ |
| if_{23} | $= \text{if } e_1 \text{ then } \mathcal{G}_o[[e_2]] (I \cup \{e_1 \leftrightarrow T\}) \text{ else } \mathcal{G}_o[[e_3]] (I \cup \{e_1 \leftrightarrow F\})$ | otherwise |
| let | $\mathcal{G}_o[[\text{let } v = e_1 \text{ in } e_2 \text{ end}]] I$ | |
| let_1 | $= \mathcal{G}_o[[\text{let } v = (\mathcal{G}_o[[e_1]]I) \text{ in } e_2 \text{ end}]] I$ | if e_1 is reducible |
| let_{1-v} | $= \mathcal{G}_o[[e_2[v := v']]] I$ | if e_1 is irreducible and e_1 is v' |
| let_{1-f} | $= \mathcal{G}_o[[\text{if } e'_1 \text{ then } (\text{let } v = e'_2 \text{ in } e_2 \text{ end}) \text{ else } (\text{let } v = e'_3 \text{ in } e_2 \text{ end})]] I$ | if e_1 is irreducible and e_1 is if e'_1 then e'_2 else e'_3 |
| let_{1-let} | $= \mathcal{G}_o[[\text{let } v' = e'_1 \text{ in } (\text{let } v = e'_2 \text{ in } e_2 \text{ end}) \text{ end}]] I$ | if e_1 is irreducible and e_1 is let $v' = e'_1$ in e'_2 end |
| let_2 | $= \mathcal{G}_o[[e_2]] I$ | if e_1 is irreducible and is not if or let , and $v \leftrightarrow_I^* e_1$ |
| let'_2 | $= \text{let } v = e_1 \text{ in } \mathcal{G}_o[[e_2]] (I \cup \{v \leftrightarrow e_1\}) \text{ end}$ | otherwise |

Figure 1: Definition of \mathcal{G}_o

i th subexpression of the top-level construct and condition it on the *irreducibility* of subexpressions 1 through $i - 1$. Operationally, when we say that subexpressions 1 through $i - 1$ are irreducible, we mean that they are the result of having already applied \mathcal{G}_o for subexpressions at those positions. We say that an expression e is *reducible* if e is not irreducible.

The meaning of v rules are straightforward. A variable that can be shown to be equal to a constant c under the assumptions of the given information set is replaced by c .

For a constructor application, \mathcal{G}_o recursively reduces each argument in turn. If a reduced argument expression is an **if**-expression or a **let**-expression, then the **if** or **let** is lifted out of the constructor.

On a primitive function application, \mathcal{G}_o first reduces each argument and, as in the case of a constructor application, lifts any **if** or **let** expressions. Then it uses the properties of the primitive function to simplify the application. We say that a primitive function application $p(e_1, \dots, e_n)$ can be *simplified (to e') under I* if we can effectively find e' such that $p(e_1, \dots, e_n) \leftrightarrow_{\mathcal{I}}^* e'$ and $t(e') \leq t(p(e_1, \dots, e_n))$. We define $\text{Simp}[\![p(e_1, \dots, e_n)]\!]I$ to be an expression e' if $p(e_1, \dots, e_n)$ can be simplified to e' under I but e' can not be further simplified under I .

On a function application, \mathcal{G}_o also first reduces each argument and, if necessary, lifts an **if** or a **let** as in the case of a constructor application. Then, it considers whether the function application needs to be unfolded. As this is the key point for \mathcal{G}_o , it will be discussed in detail below.

For an **if**-expression, \mathcal{G}_o first reduces the Boolean expression. If the reduced Boolean expression is another **if** or **let**-expression, then this **if** or **let** is lifted out of the original top-level **if**-expression. Suppose the Boolean expression is reduced already and is not an **if** or a **let**-expression. Then if the Boolean expression can be proved to be true (false) under the assumptions of the given information set, the **if** structure simplifies to the true (false) branch alone, which is then reduced. Otherwise, if the Boolean expression can not be proved to be true or false, the Boolean expression is left as a residual and \mathcal{G}_o reduces both the true and false branches. In this case, the assumption that the Boolean expression equals true (false) is added into the information set for the true (false) branch of the **if**-expression.

For a **let**-expression, \mathcal{G}_o is similar to the case of an **if**-expression. It first reduces the binding expression and, if necessary, lifts **if** and **let**-expressions. If the binding expression is reduced to a variable, then a minor simplification is made by a variable substitution. Suppose the binding expression is reduced already and is not an **if**-expression, a **let**-expression, or a variable. Then if the variable v can be proved to be equal to the binding expression, then the **let** structure simplifies to its body alone, which is then reduced. Otherwise, the binding is left as a residual and \mathcal{G}_o reduces the body of the **let**-expression with the information set enlarged by the binding. Recall that we assume

that all variables in **let** bindings are distinct, which guarantees the consistency of information sets.

We now return to the question of which function applications to unfold. Recall that we are interested in unfolding function applications that depend on both the old input x and the change parameter y .

Given an expression e , a single change operation \oplus_i , and an information set I , we say that e is *favorably decomposable (under I)* into $e_x \oplus_i e_y$ if e depends on both x and y and we can find two tuples of expressions, e_x and e_y , such that

- 1) e_x depends on x but not on y ,
- 2) e_y depends on y and possibly on x , and
- 3) $e \leftrightarrow_I^* e_x \oplus_i e_y$.

As a trivial but important example, if $e \equiv e_1 \oplus_i e_2$, where e_1 depends on x but not on y and e_2 depends on y , then e is favorably decomposable into $e_1 \oplus_i e_2$ under any information set I .

Given the input change operation $\oplus \equiv \langle \oplus_1, \dots, \oplus_n \rangle$, we say that an expression e is *most favorably decomposable* into $e_x \oplus_i e_y$ if e is favorably decomposable into $e_x \oplus_i e_y$ for some i and $t(e_x \oplus_i e_y) \leq t(e'_x \oplus_j e'_y)$ for all favorable decompositions $e'_x \oplus_j e'_y$ of e for $j = 1, \dots, n$.

In order to avoid unlimited unfoldings, we associate an *argument pattern* with each function application and unfold applications of a function at most once for each of its *favorable* argument patterns. Consider a function application $f(e_1, \dots, e_n)$. Suppose that e_1, \dots, e_n have been reduced and are not **if** or **let**-expressions. We say that $f(e_1, \dots, e_n)$ has a *favorable argument pattern* if at least one of the e_i 's either has a most favorable decomposition or depends on neither x nor y . Formally, we associate an *argument pattern* $\langle pat_1, \dots, pat_n \rangle$ with the function application $f(e_1, \dots, e_n)$, where

$$pat_i = \begin{cases} [j] & \text{if } e_i \text{ is most favorably decomposable into } e_{ix} \oplus_j e_{iy} \text{ for some } j \\ e_i & \text{if } e_i \text{ depends on neither } x \text{ nor } y \\ [0] & \text{otherwise} \end{cases}$$

We say that the argument pattern of a function application is *favorable* if not all of its components are $[0]$.

For each function f , we keep a list l_f of the argument patterns for which f has already been unfolded. The unfolding rule f_e of \mathcal{G}_o unfolds an application of f if and only if its argument pattern is favorable and does not occur in the list l_f . Note that this is merely one of many possible unfolding strategies; a discussion of unfolding strategies appears in Section 7.

The unfolding rule is the only possible source of nontermination. We can guarantee that \mathcal{G}_o always terminates by forbidding the unfolding of function applications caused by arguments that depend on neither x nor y , i.e., we could define pat_i to be $[j]$ if e_i is most favorably decomposable

into $e_{ix} \oplus_j e_{iy}$ and $[0]$ otherwise. For this notion of favorable argument pattern, each function has only a fixed number of favorable argument patterns, so termination is assured. This choice omits the class of straightforward optimization — evaluation by unfolding and simplification. Of course, allowing such unfoldings coupled with some heuristics to prevent divergence could be an effective implementation strategy.

We can prove the following proposition by an induction on the structure of expressions, following the \mathcal{G}_o rules, and making use of the fact that the underlying logic is compatible with the evaluation semantics of our language.

Proposition 5.2 *For any expression e and information set I of e , let z_1, \dots, z_k be all the free variables in e and e_I . If $\mathcal{G}_o[[e]I]$ terminates, then, $\forall v_z$, if $\mathcal{E}[[e_I]] [z \mapsto v_z] = T$, then $\forall v$, $\mathcal{E}[[e]] [z \mapsto v_z] \Downarrow = v \implies \mathcal{E}[[\mathcal{G}_o[[e]I]] [z \mapsto v_z] \Downarrow = v$ and $t(\mathcal{G}_o[[e]I]) \leq t(e)$.*

Therefore, if the application of \mathcal{G}_o on $f_0(x \oplus y)$ terminates and we get $e^\circ \equiv \mathcal{G}_o[[f_0(x \oplus y)]] \emptyset$, then we have the following corollary.

Corollary 5.3 *If $\mathcal{G}_o[[f_0(x \oplus y)]] \emptyset$ terminates with $e^\circ \equiv \mathcal{G}_o[[f_0(x \oplus y)]] \emptyset$, then we have, $\forall v_x, v_y, v$,*

$$\mathcal{E}[[f_0(x \oplus y)]] [x \mapsto v_x, y \mapsto v_y] \Downarrow = v \implies \mathcal{E}[[e^\circ]] [x \mapsto v_x, y \mapsto v_y] \Downarrow = v$$

and $t(e^\circ) \leq t(f_0(x \oplus y))$.

5.4 Step 2: Introducing the cached-result parameter r

Step 2 defines the function f_0^* such that $f_0^*(x, y, r)$ computes $f_0(x \oplus y)$ when r is the cached result $f_0(x)$. The extra parameter r is introduced as an explicit name for the cached result $f_0(x)$. Our motivation for introducing r is so that, in a later step, some computations performed by $f_0(x \oplus y)$ may be replaced by expressions involving the cached result r , provided that such expressions compute the same values and are at least as fast as doing the original computations.

We define

$$f_0^*(x, y, r) = e^*, \tag{13}$$

where e^* is a transformed version of e° . The transformation of e° involves replacing some calls $f_0(e)$ with recursive calls $f_0^*(e_x, e_y, f_0(e_x))$, where e_x and e_y are a *definition candidate pair* for e , defined below, such that $e_x \oplus e_y$ is equivalent to e . Note that this transformation preserves the invariant that the third argument to f_0^* is always f_0 applied to the first argument of f_0^* . Also note that, in

the body of f_0^* as defined in this step, parameter r is not used in any substantive way; it will be used by Step 3 to replace, in the body of f_0^* , some computations on x by expressions involving r .

The transformation of e^o into e^* is implemented by \mathcal{D}_t , which has type $Exp \rightarrow Info \rightarrow Exp$. Expression e^* , the body of function f_0^* , is obtained by applying \mathcal{D}_t to e^o and the null information set \emptyset :

$$e^* \equiv \mathcal{D}_t[[e^o]] \emptyset. \quad (14)$$

Figure 2 defines \mathcal{D}_t . The names of rules have the same meaning as for the \mathcal{G}_o transformation, with the addition that the subscript s indicates substitution of a function application.

| Name | Transformation | Condition |
|-------|--|--|
| v | $\mathcal{D}_t[[v]] I = v$ | |
| c | $\mathcal{D}_t[[c(e_1, \dots, e_n)]] I = c(\mathcal{D}_t[[e_1]]I, \dots, \mathcal{D}_t[[e_n]]I)$ | |
| p | $\mathcal{D}_t[[p(e_1, \dots, e_n)]] I = p(\mathcal{D}_t[[e_1]]I, \dots, \mathcal{D}_t[[e_n]]I)$ | |
| f_s | $\mathcal{D}_t[[f(e_1, \dots, e_n)]] I = f_0^*(e_x, e_y, f(e_x))$ | if $f \equiv f_0$ and $\langle e_x, e_y \rangle$ is a definition candidate pair for $f(e_1, \dots, e_n)$ |
| f_n | $\mathcal{D}_t[[f(e_1, \dots, e_n)]] I = f(\mathcal{D}_t[[e_1]]I, \dots, \mathcal{D}_t[[e_n]]I)$ | otherwise |
| if | $\mathcal{D}_t[[if\ e_1\ then\ e_2\ else\ e_3]] I = if\ \mathcal{D}_t[[e_1]]I\ then\ \mathcal{D}_t[[e_2]]\ (I \cup \{e_1 \leftrightarrow T\})\ else\ \mathcal{D}_t[[e_3]]\ (I \cup \{e_1 \leftrightarrow F\})$ | |
| let | $\mathcal{D}_t[[let\ v = e_1\ in\ e_2\ end]] I = let\ v = \mathcal{D}_t[[e_1]]I\ in\ \mathcal{D}_t[[e_2]]\ (I \cup \{v \leftrightarrow e_1\})\ end$ | |

Figure 2: Definition of \mathcal{D}_t

All rules other than those dealing with function applications are straightforward. They just collect and pass on information sets for the expressions being transformed. If there is no function call to f_0 in e^o , then $e^* \equiv e^o$.

Consider a (recursive) function application $f_0(e_1, \dots, e_n)$ in the expression e^o with the collected information I . Let $e \equiv \langle e_1, \dots, e_n \rangle$. If there exist a pair of tuples e_x and e_y such that e_i is favorably decomposable into $e_x \oplus_i e_y$ under I for each $i = 1, \dots, n$ and $t(e_x \oplus e_y) \leq t(e)$, then we call $\langle e_x, e_y \rangle$ a *definition candidate pair* for $f_0(e)$, and replace the occurrence of $f_0(e)$ with $f_0^*(e_x, e_y, f_0(e_x))$.

Function f_0^* is not necessarily a differentiated version of f_0 because, although transformation \mathcal{D}_t preserves semantics, it may actually result in a slower program. To compare the time of e^* and e^o , consider a function application $f_0^*(e_x, e_y, f_0(e_x))$ in e^* that replaces function application $f_0(e)$ in e^o . The fact that $\langle e_x, e_y \rangle$ is a definition candidate pair for $f_0(e)$ guarantees that computing e_x and e_y is at least as fast as computing e . However, $f_0(e_x)$ may need a substantial amount of time to compute and hence transformation \mathcal{D}_t produces a slower program. This apparent step

backwards is remedied in Step 3, which restores $f_0^*(e_x, e_y, f_0(e_x))$ to $f_0(e)$ unless it can be replaced by $f_0^*(e_x, e_y, e_r)$, where e_r computes $f_0(e_x)$ rapidly using the cached value r and $f_0^*(e_x, e_y, f_0(e_x))$ can compute at least as fast as $f_0(e)$ by making use of e_r .

Note that, among various heuristics that can be used to find a definition candidate pair for $f(e)$ under information I , there are two special cases that are most useful:

1. If there exists a right-reverse operation \ominus of \oplus such that $y = z \ominus x$ iff $x \oplus y = z$, and $e \ominus x$ can be simplified element-wise under I to form e_y such that $t(x \oplus e_y) \leq t(e)$, then $\langle x, e_y \rangle$ is a definition candidate pair for $f_0(e)$. In this case, transformation \mathcal{D}_t may replace $f_0(e)$ by $f_0^*(x, e_y, f_0(x))$, and Step 3 can subsequently replace $f_0(x)$ by the cached result r .
2. If there exists a left-reverse operation \ominus of \oplus such that $x = z \ominus y$ iff $x \oplus y = z$, and $e \ominus y$ can be simplified element-wise under I to form e_x such that e_x does not depend on y and $t(e_x \oplus y) \leq t(e)$, then $\langle e_x, y \rangle$ is a definition candidate pair for $f_0(e)$. In this case, transformation \mathcal{D}_t may replace $f_0(e)$ by $f_0^*(e_x, y, f_0(e_x))$.

For the transformation \mathcal{D}_t described above, we have the following proposition:

Proposition 5.4 *With the function definition $f_0^*(x, y, r) = e^*$, we have, $\forall v_x, v_y, v_1, v_2$,*

$$\mathcal{E}[[e^0]] [x \mapsto v_x, y \mapsto v_y] \downarrow = v_1 \wedge \mathcal{E}[[e^*]] [x \mapsto v_x, y \mapsto v_y] \downarrow = v_2 \implies v_1 = v_2.$$

5.5 Cache sets

We would like to find, in the computation of $f_0(x \oplus y)$ performed by $f_0^*(x, y, r) \equiv e^*$, expressions whose values can be retrieved quickly from the cached result r . Of course, literal occurrences of the subexpression $f_0(x)$ in e^* can be replaced by r . But we aim to locate many other opportunities to use the cached result r . For example,

- Expressions that are provably equal to but not identical to $f_0(x)$ can be replaced by r .
- If $f_0(x)$ is provably equal to $g(h(x))$ and g has an inverse g^{-1} such that $t(g^{-1}(r)) \leq t(h(x))$ for $f_0(x) = r$, then any expressions provably equal to $h(x)$ can be replaced by $g^{-1}(r)$.

Such equalities need not hold globally throughout e^* . Rather, they only need to hold in the local context of the expression that is a candidate for replacement.

For e being a subexpression of e^* , let I_e be the information set at e , say as computed by Step 1 or 2. We associate a *cache set* C_e with e . We use *Cache* to denote the domain of cache sets. For any $C_e \in \text{Cache}$, C_e is a set of expression pairs $\langle e', r' \rangle$ such that,

- 1) $e' \leftrightarrow_{I_e}^* r'$, and
- 2) r is the only free variable of r'

In general, the expression pairs in C_e may have no relation to e other than their validity with respect to I_e . Note that, although all pairs $\langle e', r' \rangle$ in C_e satisfy 1) and 2) above, C_e is not necessarily a complete set of such pairs. The computation of C_e will be given in Step 3, where a cache set is collected for the expression e and is then used to determine whether e can be replaced by a use of the cached result r .

For any subexpression e of e^* , r' is called a *substitution candidate* for e iff

- 1) $\exists e', e \leftrightarrow_I^* e'$ and $\langle e', r' \rangle \in C_e$, and
- 2) $t(r') \leq t(e)$.

We use $e \rightarrow_{IC}^* r'$ to denote that r' is a substitution candidate for e and $t(r') \leq t(r'')$ for all substitution candidates r'' of e that we can find.

5.6 Step 3: Using the Cache

We are now ready to redefine function f_0^* so that it makes use of the cached-result parameter r . Step 2 defined f_0^* by $f_0^*(x, y, r) = e^*$. We now redefine

$$f_0^*(x, y, r) = e^t, \quad (15)$$

where e^t is derived from e^* by transformation \mathcal{S}_t :

$$e^t \equiv \mathcal{S}_t[[e^*]] \emptyset \{\langle f_0(x), r \rangle\}, \quad (16)$$

and \mathcal{S}_t has type $Exp \rightarrow Info \rightarrow Cache \rightarrow Exp$. For each subexpression e of e^* , \mathcal{S}_t collects the information set I_e from the context of e , computes the cache set C_e , and, where possible, substitutes e with its fastest substitution candidate.

Figure 3 defines \mathcal{S}_t . The first column gives a name for each rule. The subscripts of rule names have the same significance as those for \mathcal{G}_o . The subscript r corresponds to a rule that does a direct substitution using the cached result. The subscript rp corresponds to a rule that does a substitution for a predicate application.

\mathcal{S}_t works inductively on the structure of expressions. For each expression e , it substitutes e with r' if $e \rightarrow_{IC}^* r'$. Otherwise, it examines e 's subexpressions. Rule v_n is the base case. Rule c_i involves a straightforward application of \mathcal{S}_t on each argument.

For a primitive function application $p(e_1, \dots, e_n)$, if p is a predicate and $p(e_1, \dots, e_n)$ is *predicate substitutable* by $p(r'_1, \dots, r'_n)$, then rule p_{rp} replaces $p(e_1, \dots, e_n)$ by $p(r'_1, \dots, r'_n)$. The formal notion

| Name | Transformation | Condition |
|-----------|--|--|
| v | $S_t \llbracket v \rrbracket I C$ | |
| v_r | $= r'$ | if $v \rightarrow_{IC}^* r'$ |
| v_n | $= v$ | otherwise |
| c | $S_t \llbracket c(e_1, \dots, e_n) \rrbracket I C$ | |
| c_r | $= r'$ | if $c(e_1, \dots, e_n) \rightarrow_{IC}^* r'$ |
| c_i | $= c(S_t \llbracket e_1 \rrbracket I C, \dots, S_t \llbracket e_n \rrbracket I C)$ | otherwise |
| p | $S_t \llbracket p(e_1, \dots, e_n) \rrbracket I C$ | |
| p_r | $= r'$ | if $p(e_1, \dots, e_n) \rightarrow_{IC}^* r'$ |
| p_{rp} | $= p(r'_1, \dots, r'_n)$ | else if $p(e_1, \dots, e_n)$ is predicate substitutable by $p(r'_1, \dots, r'_n)$ under I and C |
| p_i | $= \text{Simp} S_t \llbracket p(S_t \llbracket e_1 \rrbracket I C, \dots, S_t \llbracket e_n \rrbracket I C) \rrbracket I C$ | otherwise |
| f | $S_t \llbracket f(e_1, \dots, e_n) \rrbracket I C$ | |
| f_r | $= r'$ | if $f(e_1, \dots, e_n) \rightarrow_{IC}^* r'$ |
| f_{rp} | $= f(r'_1, \dots, r'_n)$ | else if $f(e_1, \dots, e_n)$ is predicate substitutable by $f(r'_1, \dots, r'_n)$ under I and C |
| f_i | $= f(S_t \llbracket e_1 \rrbracket I C, \dots, S_t \llbracket e_n \rrbracket I C)$ | else if $f \neq f_0^*$ or $f \equiv f_0^*$ and $f_0^*(e_1, \dots, e_n)$ is preferred over the original function application $f_0(e)$ |
| f'_i | $= S_t \llbracket f_0(e) \rrbracket I C$ | otherwise |
| if | $S_t \llbracket \text{if } e_1 \text{ then } e_2 \text{ else } e_3 \rrbracket I C$ | |
| if_r | $= r'$ | if $\text{if } e_1 \text{ then } e_2 \text{ else } e_3 \rightarrow_{IC}^* r'$ |
| if_{23} | $= \text{if } e'_1 \text{ then } S_t \llbracket e_2 \rrbracket I_2 C_2 \text{ else } S_t \llbracket e_3 \rrbracket I_3 C_3$ | otherwise |
| | where $e'_1 = S_t \llbracket e_1 \rrbracket I C$, $I_2 = I \cup \{e_1 \leftrightarrow T\}$, $C_2 = C(C, I_2)$, $I_3 = I \cup \{e_1 \leftrightarrow F\}$, $C_3 = C(C, I_3)$ | |
| let | $S_t \llbracket \text{let } v = e_1 \text{ in } e_2 \text{ end} \rrbracket I C$ | |
| let_r | $= r'$ | if $\text{let } v = e_1 \text{ in } e_2 \text{ end} \rightarrow_{IC}^* r'$ |
| let'_2 | $= \text{let } v = e'_1 \text{ in } S_t \llbracket e_2 \rrbracket I_2 C_2 \text{ end}$ | otherwise |
| | where $e'_1 = S_t \llbracket e_1 \rrbracket I C$, $I_2 = I \cup \{v \leftrightarrow e_1\}$, $C_2 = C(C, I_2)$ | |
| $simp$ | $\text{Simp} S_t \llbracket p(e_1, \dots, e_n) \rrbracket I C$ | |
| $simp_s$ | $= \text{Simp} p \llbracket p(e_1, \dots, e_n) \rrbracket I C$ | if $p(e_1, \dots, e_n)$ can be simplified with substitution under I and C |
| $simp_n$ | $= e$ | otherwise |

Figure 3: Definition of S_t

of predicate substitution will be explained later, but the idea is that $p(e_1, \dots, e_n)$ may be replaced by $p(r'_1, \dots, r'_n)$ even if each e_i may not be replaced by r'_i . For example, if sorting a list x returns a sorted list r , i.e., $\text{sort}(x) = r$, then although, in general, x is not equivalent to r , $\text{null}(x)$ is equivalent to $\text{null}(r)$ and thus can be replaced by $\text{null}(r)$. If none of the rules p_r or p_{rp} applies, then p_i applies \mathcal{S}_t straightforwardly on each argument. $\text{Simp}\mathcal{S}_t$ simplifies the resulting primitive function application by making use of cache C and simplifier Simp described in Step 1. It is defined in rule simp . Given information set I and cache set C , let I^C denote the set $I \cup \{e \leftrightarrow r \mid (e, r) \in C\}$. We say that primitive function application $p(e_1, \dots, e_n)$ can be *simplified by substitution under I and C* if $p(e_1, \dots, e_n)$ can be simplified under I^C ; we define $\text{Simp}\mathcal{S}_t[[p(e_1, \dots, e_n)]] I C = e'$ if $\text{Simp}[[p(e_1, \dots, e_n)]] I^C = e'$.

For a function application $f(e_1, \dots, e_n)$, if f is a Boolean function and $f \not\equiv f_0^*$, then f_{rp} also tries a predicate substitution as rule p_{rp} does. If $f \equiv f_0^*$, then, as a result of Step 2, $f(e_1, \dots, e_n) \equiv f_0^*(e_x, e_y, f_0(e_x))$ for some e_x and e_y . In this case, \mathcal{S}_t makes sure that the function application of f_0^* that eventually appears in e^t can be computed at least as fast as the original function application of f_0 . Formally, if $e_r = \mathcal{S}_t[[f_0(e_x)]]IC$, we say that $f_0^*(e_x, e_y, f_0(e_x))$ is *preferred* over the original function application $f_0(e)$ if $t(f_0^*(e_x, e_y, e_r)) \leq t(f_0(e))$. Testing for such preference requires us to compare the definitions of f_0^* and f_0 ; Proposition 5.5 gives a heuristic that can always be applied for this purpose.

Proposition 5.5 *For e^o as obtained from (12), if $t(e^o) \leq t(f_0(x \oplus y))$ and $t(e_r) \leq t(e)$ for each $f_0(e)$ in e^o that is replaced by $f_0^*(e_x, e_y, e_r)$ in e^t , then with $f_0^*(x, y, r) = e^t$, we have $t(e^t) \leq t(e^o)$ and $t(f_0^*(e_x, e_y, e_r)) \leq t(f_0(e))$.*

This heuristic requires that the time to compute e^o be no greater than the time to compute $f_0(x \oplus y)$ in an absolute sense, i.e., a slowdown by a constant factor is not allowed. Since the possible slowdown is caused by computing argument expressions more than once after function unfoldings, we only need to guarantee that these expressions can be simplified in context so that together they take at most the time of computing the argument expressions once. If $f_0^*(e_x, e_y, f_0(e_x))$ is *preferred* over the original function application $f_0(e)$, then the function application of f_0^* is retained in the resulting expression e^t and \mathcal{S}_t is applied to each argument; otherwise, the original function application $f_0(e)$ is restored and \mathcal{S}_t is applied to it. Note that function applications are not further unfolded by \mathcal{S}_t , as all necessary unfoldings have been done by \mathcal{G}_o in Step 1.

For **if** and **let**-expressions, \mathcal{S}_t applies similar ideas. It collects information from the logical structure of the expression in the set I , enlarges the cache C using a closure-like function \mathcal{C} , defined

below, and continues making substitutions in the two branches of an **if**-expression or the body of a **let**-expression, respectively.

In order to define the notion of predicate substitution and the function \mathcal{C} , we need the aGPE transformation \mathcal{G}_a as a subroutine. Given an expression e' and an information set I , $\mathcal{G}_a[[e']]I$ specializes e' with respect to I . The basic idea of using \mathcal{G}_a is this: given an expression e in the context of information I , if there is $\langle e', r' \rangle \in C$ such that $\mathcal{G}_a[[e']]I = e$, then e can be replaced by r' provided that $t(r') \leq t(e)$.

Given a primitive function application or a function application $P(e_1, \dots, e_n)$ for some predicate P , let I and C be the information set and cache set of $P(e_1, \dots, e_n)$. We say that $P(e_1, \dots, e_n)$ is *predicate substitutable* with $P(r'_1, \dots, r'_n)$ under I and C if there are $\langle e'_i, r'_i \rangle \in C$ such that

$$\begin{aligned}
& 1) \mathcal{G}_a[[P(e'_1, \dots, e'_n)]](I \cup \{P(e_1, \dots, e_n) \leftrightarrow T\}) = T, \\
& \quad \mathcal{G}_a[[P(e'_1, \dots, e'_n)]](I \cup \{P(e_1, \dots, e_n) \leftrightarrow F\}) = F, \\
& 2) $P(e'_1, \dots, e'_n)$ is defined if and only if $P(e_1, \dots, e_n)$ is defined, and \\
& 3) $t(P(r'_1, \dots, r'_n)) \leq t(P(e_1, \dots, e_n))$.
\end{aligned} \tag{17}$$

In the matrix multiplication example of Section 3, consider $null(C)$ in expression (9). We have $\langle mtxMul(C, R), r \rangle$ in the cache set, and

$$\mathcal{G}_a[[null(mtxMul(C, R))]] \{null(C) \leftrightarrow T\} = T, \quad \mathcal{G}_a[[null(mtxMul(C, R))]] \{null(C) \leftrightarrow F\} = F,$$

since when $null(C)$ equals T , $mtxMul(C, R)$ is specialized to nil and thus $null(mtxMul(C, R))$ to T , and when $null(C)$ equals F , $mtxMul(C, R)$ is specialized to a *cons* structure and thus $null(mtxMul(C, R))$ to F . Therefore, the expression $null(C)$ can actually be replaced by $null(r)$.

The function \mathcal{C} used in **if** and **let** rules has type $Cache \rightarrow Info \rightarrow Cache$. For any $C \in Cache, I \in Info$, we call $\mathcal{C}(C, I)$ the *closure of C under I* and define $\mathcal{C}(C, I) = C \cup C' \cup C''$, where

$$\begin{aligned}
C' &= \{\langle \mathcal{G}_a[[e']]I, r' \rangle \mid \langle e', r' \rangle \in C\}, \\
C'' &= \{\langle e'_1, g^{-1}(r') \rangle \mid \langle e', r' \rangle \in C, \mathcal{G}_a[[e']]I \leftrightarrow_I^* g(e'_1)\}.
\end{aligned} \tag{18}$$

For any $\langle e', r' \rangle \in C$, if \mathcal{G}_a specializes e' to e'' under I , then $\langle e'', r' \rangle$ is added into C' . If e'' is provably equal to $g(e'_1)$ for some primitive function g and g has an inverse g^{-1} , then $\langle e'_1, g^{-1}(r') \rangle$ is added into C'' . In particular, if $e'' \equiv c(e'_1, \dots, e'_n)$ for some constructor c and c 's corresponding destructors are c_i^{-1} 's, then $\langle e'_i, c_i^{-1}(r') \rangle$ is added into C'' for $i = 1, \dots, n$. In the matrix multiplication example of Section 3, the expression in the false branch of expression (6) has $I = \{null(C) \leftrightarrow F\}$ and $C = \{\langle mtxMul(C, R), r \rangle\}$. Then $C' = \{\langle cons(rowMul(car(C), R), mtxMul(cdr(C), R)) \rangle\}$ by (7)

and $C'' = \{\langle \text{rowMul}(\text{car}(C), R), \text{car}(r) \rangle, \langle \text{mtxMul}(\text{cdr}(C), R), \text{cdr}(r) \rangle\}$ by (8). Thus, the cache set associated with that expression would be $C \cup C' \cup C''$ by (18).

From the above definitions, we see how \mathcal{G}_a transformations are used. The goal is to help find more expressions whose values can be retrieved from a cached result. Therefore, an unfolding by \mathcal{G}_a is beneficial for this purpose only if the unfolded expression can be specialized to a new expression under the given information I . In general, the new expression can be obtained in one of the following two ways:

- 1) simplification of primitive function applications using *Simp* as rule p_e in \mathcal{G}_o ;
- 2) specialization of **if** or **let** expressions using rules like if_2 , if_3 , or let_2 in \mathcal{G}_o .

Therefore, we unfold a function application only if the unfolded expression can be simplified using *Simp* or specialized to a strict subexpression of itself. Other than the unfolding rule, \mathcal{G}_a does most of the other transformations and simplifications in a fashion similar to \mathcal{G}_o .

Figure 4 defines \mathcal{G}_a , which has type $Exp \rightarrow Info \rightarrow Exp$. The definition of \mathcal{G}_a is mostly the same as for \mathcal{G}_o ; it differs only in its treatment of function applications and **let**-expressions. The notion of *irreducibility* for \mathcal{G}_a is similar to irreducibility for \mathcal{G}_o .

| Name | Transformation | Condition |
|-------------|---|---|
| f | $\mathcal{G}_a[[f(e_1, \dots, e_n)]] I$ | |
| f_e | $= \mathcal{G}_a[[e[x_1 := e_1, \dots, x_n := e_n]]] I$ where f is defined by $f(x_1, \dots, x_n) = e$ | if e_1, \dots, e_n are irreducible, not if or let , but $f(e_1, \dots, e_n)$ is expandable. |
| let | $\mathcal{G}_a[[\text{let } v = e_1 \text{ in } e_2 \text{ end}]] I$ | |
| let_{1-v} | < this case is deleted > | |
| let'_2 | $= \mathcal{G}_a[[e_2[v := e_1]]] I$ | otherwise |
| * | all other rules are the same as those for \mathcal{G}_o . | |

Figure 4: Definition of \mathcal{G}_a

Consider a function application $f(e_1, \dots, e_n)$ in the context of information set I , where the e_i 's are already reduced by \mathcal{G}_a and are not **if** or **let**-expressions. Suppose e is the expression obtained by unfolding f on $\langle e_1, \dots, e_n \rangle$. We say that $f(e_1, \dots, e_n)$ is *expandable by \mathcal{G}_a* if e can be simplified by *Simp* under I or specialized to be a strict subexpression of itself under I . \mathcal{G}_a unfolds a function application $f(e_1, \dots, e_n)$ if and only if it is expandable. Recall that \mathcal{G}_o unfolds $f(e_1, \dots, e_n)$ if and only if its argument pattern is favorable and occurs the first time for f . Here again, the unfolding rule f_e is a potential source of nontermination.

For a **let**-expression, \mathcal{G}_a applies rules $let_1, let_{1-if}, let_{1-let}, let_2$ in turn, if possible. At the end, instead of leaving the binding in the residual expression as \mathcal{G}_o does, \mathcal{G}_a replaces occurrences of each

defined variable v in e_2 by its defining expression e_1 . This is merely an implementation strategy so that, when using the cache set, we do not have to keep track of equations introduced by the bindings of **let**-expressions. When a program has many **let**-expressions, this may cause code blow-up. However, when a program has a lot of **let**-expressions, we expect that many of them can be eliminated by the let_2 rule in the context of the information collected for the \mathcal{G}_a transformation.

We can prove the following proposition by induction on the structure of expressions:

Proposition 5.6 *For any expression e and information set I of e , let z_1, \dots, z_k be all the free variables in e and e_I . If $\mathcal{G}_a[[e]]I$ terminates, then we have $\forall w$, if $\mathcal{E}[[e_I]] [z \mapsto v_z] = T$, then $\forall v$, $\mathcal{E}[[e]] [z \mapsto v_z] \downarrow = v \implies \mathcal{E}[[\mathcal{G}_a[[e]]I]] [z \mapsto v_z] \downarrow = v$.*

Corollary 5.7 *For any cache set C and information set I , let z_1, \dots, z_k be all the free variables in e_I and e' and r' 's for all $\langle e', r' \rangle \in C$. Then we have $\forall v_z$, if $\mathcal{E}[[e_I]] [z \mapsto v_z] = T$ and $\forall \langle e', r' \rangle \in C$, $\mathcal{E}[[e']] [z \mapsto v_z] \downarrow = \mathcal{E}[[r']] [z \mapsto v_z] \downarrow$, then $\forall \langle e', r' \rangle \in C(C, I)$, $\mathcal{E}[[e']] [z \mapsto v_z] \downarrow = \mathcal{E}[[r']] [z \mapsto v_z] \downarrow$.*

We redefine f_0^* to be $f_0^*(x, y, r) = e^t$. We have the following proposition for \mathcal{S}_t transformation:

Proposition 5.8 *For any expression e , information set I and cache set C of e , let z_1, \dots, z_k be all the free variables in e , e_I , and e' and r' 's for all $\langle e', r' \rangle \in C$. Let $e_1 \equiv \mathcal{D}_t[[e]]I$. If $\mathcal{S}_t[[e_1]]IC$ terminates with $e_2 \equiv \mathcal{S}_t[[e_1]]IC$, then with the definition $f_0^*(x, y, r) = e^t$, we have $\forall v_z$, if $\mathcal{E}[[e_I]] [z \mapsto v_z] = T$ and $\forall \langle e', r' \rangle \in C$, $\mathcal{E}[[e']] [z \mapsto v_z] \downarrow = \mathcal{E}[[r']] [z \mapsto v_z] \downarrow$, then $\forall v_1, v_2$, $\mathcal{E}[[e]] [z \mapsto v_z] \downarrow = v_1 \wedge \mathcal{E}[[e_2]] [z \mapsto v_z] \downarrow = v_2 \implies v_1 = v_2$ and $t(e_2) \leq t(e)$.*

Therefore, if the \mathcal{G}_a transformations issued by \mathcal{S}_t terminate and we get $e^t \equiv \mathcal{S}_t[[e^*]] \emptyset \{\langle f_0(x), r \rangle\}$, then we have the following corollary:

Corollary 5.9 *If $\mathcal{S}_t[[e^*]] \emptyset \{\langle f_0(x), r \rangle\}$ terminates with $e^t \equiv \mathcal{S}_t[[e^*]] \emptyset \{\langle f_0(x), r \rangle\}$, then with the definition $f_0^*(x, y, r) = e^t$, we have $\forall v_x, v_y, v_r, v_1, v_2$, if $\mathcal{E}[[f_0(x)]] [x \mapsto v_x] \downarrow = v_r$, then*

$$\mathcal{E}[[e^o]] [x \mapsto v_x, y \mapsto v_y] \downarrow = v_1 \wedge \mathcal{E}[[e^t]] [x \mapsto v_x, y \mapsto v_y, r \mapsto v_r] \downarrow = v_2 \implies v_1 = v_2.$$

and $t(e^t) \leq t(e^o)$.

Finally, we have our main result, which follows directly from Corollary 5.3 and Corollary 5.9: the semantics of the programs are preserved by the series of transformations above and computations using the differentiated versions are at least as fast as computations using the original programs.

Theorem 5.10 *If the oGPE and aGPE transformations terminate, then with the function definition $f_0^*(x, y, r) = e^t$, we have,*

- (1) $\forall v_x, v_y, v_r$, if $\mathcal{E}[[f_0(x)]] [x \mapsto v_x] \downarrow = v_r$, then $\forall v_1, v_2$,
 $\mathcal{E}[[f_0(x \oplus y)]] [x \mapsto v_x, y \mapsto v_y] \downarrow = v_1 \wedge \mathcal{E}[[f_0^*(x, y, r)]] [x \mapsto v_x, y \mapsto v_y, r \mapsto v_r] \downarrow = v_2 \implies v_1 = v_2$,
- (2) $t(f_0^*(x, y, r)) \leq t(f_0(x \oplus y))$.

5.7 Step 4: Redundant Code Elimination

The purpose of this last step is to optimize the function $f_0^*(x, y, r) = e^t$ in the context of the set of functions F . The optimizations mainly involve redundant code elimination. At least the following three kinds of code may become redundant as a result of the above transformations.

First, some function definitions in F may have become irrelevant to the evaluation of $f_0^*(x, y, r)$. This is because: After $f_0(x \oplus y)$ is expanded into e^o by Step 1, f_0^* with r is introduced in Step 2 to form e^* and then used in Step 3 to replace some expressions in e^* to form e^t . The replaced expressions may be function applications. If all calls to a function are replaced, then the definition of the function becomes redundant. In particular, some function calls to the function f_0 are possibly replaced in Step 2 by calls to f_0^* . As a result, f_0 may become redundant even without substitutions in Step 3.

Second, some computations inside remaining functions may also have become redundant. This is obvious. Since some expressions are replaced by calls to f_0^* or expressions involving r , some other expressions whose results are used for the computations of these replaced expressions may simply become unnecessary.

Third, some parameters of remaining functions may also have become redundant. This occurs because some computations dependent on the parameter x are altered by the replacements in Step 3 to be dependent on the cache parameter r . Hence all or part of the parameter x may have become irrelevant to the final computation.

There has been considerable work on such optimizations [1]. Standard data-flow and control-flow analysis algorithms can be adapted to reason about the dependencies within the programs and recognize the redundant code to be eliminated.

After optimization, we get the resulting function definition $f'_0(\bar{x}, y, r) = e'$, where $\bar{x} = \langle x_{i_1}, \dots, x_{i_k} \rangle$, $1 \leq i_1 < \dots < i_k \leq n$, in the context of a set, say F' (including f'_0), of function definitions. f'_0 is a differentiated version of f_0 . If the correctness of these standard optimizations guarantees that $\forall v_x, v_y, v_r, v$,

$$\mathcal{E}[[f_0^*(x, y, r)]] [x \mapsto v_x, y \mapsto v_y, r \mapsto v_r] \downarrow = v \iff \mathcal{E}[[f'_0(\bar{x}, y, r)]] [x \mapsto v_x, y \mapsto v_y, r \mapsto v_r] \downarrow = v,$$

and $t(f'_0(\bar{x}, y, r)) \leq t(f_0^*(x, y, r))$, then combining this with Theorem 5.10, we have $\forall v_x, v_y, v_r$, if

$\mathcal{E}[[f_0(x)]] [x \mapsto v_x] \Downarrow = v_r$, then $\forall v_1, v_2$,

$$\mathcal{E}[[f_0(x \oplus y)]] [x \mapsto v_x, y \mapsto v_y] \Downarrow = v_1 \wedge \mathcal{E}[[f'_0(\bar{x}, y, r)]] [x \mapsto v_x, y \mapsto v_y, r \mapsto v_r] \Downarrow = v_2 \implies v_1 = v_2,$$

and $t(f'_0(\bar{x}, y, r)) \leq t(f_0(x \oplus y))$.

Note that the optimizations of this last step try to make the resulting function f'_0 depend as little on the input x as possible by eliminating those x_i 's with $i \in \{1, \dots, n\} - \{i_1, \dots, i_k\}$. If x can be entirely eliminated, then $f'_0(y, r)$ is a complete differentiated version of f_0 .

6 Example: Derivation of Insert from Sort

Suppose we are given a program for select sort consisting of the function definitions for *sort*, *least*, and *rest*. *Sort* takes a list x of numbers and returns a sorted list r of these numbers. Let the change in the input of *sort* be that an extra element i is added at the beginning of the list x . Formally, we have

function definition:

```

sort(x) =  if null(x) then nil
           else let k = least(x)
                in cons(k, sort(rest(x, k))) end
least(x) = if null(cdr(x)) then car(x)
           else let s = least(cdr(x))
                in if car(x) < s then car(x) else s end
rest(x, k) = if k = car(x) then cdr(x)
             else cons(car(x), rest(cdr(x), k))

```

cached result:

```
sort(x) = r
```

change in input:

```
x' = cons(i, x)
```

We apply the series of transformations discussed in Section 5 to derive a differentiated version of *sort*. Instead of sorting the new input $\text{cons}(i, x)$ from scratch, the derived program actually inserts the element i into the sorted list r . Therefore, it is a complete differentiated version of *sort*. In fact, the derived program is also a complete incremental version of *sort*.

Step 1. Applying \mathcal{G}_o on $\text{sort}(\text{cons}(i, x))$ produces the following transformations. The most substantial steps are marked with *.

$$\begin{aligned}
& \mathcal{G}_o[[\text{sort}(\text{cons}(i, x))]]\emptyset \\
&= \mathcal{G}_o[[\text{if null}(\text{cons}(i, x)) \text{ then nil} \\
&\quad \text{else let } k = \text{least}(\text{cons}(i, x)) \\
&\quad\quad \text{in cons}(k, \text{sort}(\text{rest}(\text{cons}(i, x), k))) \\
&\quad\quad \text{end}]]\emptyset & (f_e)^* \\
&= \mathcal{G}_o[[\text{let } k = \text{least}(\text{cons}(i, x)) \\
&\quad \text{in cons}(k, \text{sort}(\text{rest}(\text{cons}(i, x), k))) \\
&\quad \text{end}]]\emptyset & (if_3)^* \\
&= \mathcal{G}_o[[\text{let } k = \mathcal{G}_o[[\text{least}(\text{cons}(i, x))]]\emptyset \\
&\quad \text{in cons}(k, \text{sort}(\text{rest}(\text{cons}(i, x), k)))
\end{aligned}$$


```

      end]]I3
    end
    where I2 = {null(x) ↔ F, s ↔ least(x), i < s ↔ T}
          I3 = {null(x) ↔ F, s ↔ least(x), i < s ↔ F}
  = if null(x) then cons(i, nil)
    else let s = least(x)
      in if i < s then Go[[cons(i, sort(rest(cons(i, x), i)))] I2
        else Go[[cons(s, sort(rest(cons(i, x), s)))] I3
      end
    = if null(x) then cons(i, nil)
      else let s = least(x)
        in if i < s then Go[[cons(i, sort(x))] I2
          else Go[[cons(s, sort(rest(cons(i, x), s)))] I3
        end
    = if null(x) then cons(i, nil)
      else let s = least(x)
        in if i < s then cons(i, sort(x))
          else Go[[cons(s, sort(Go[[rest(cons(i, x), s)] I3)] I3)] I3
        end
    = if null(x) then cons(i, nil)
      else let s = least(x)
        in if i < s then cons(i, sort(x))
          else Go[[cons(s, sort(Go[[if s = car(cons(i, x))
            then cdr(cons(i, x))
            else cons(car(cons(i, x)),
              rest(cdr(cons(i, x), s))
            ] I3)] I3)] I3
        end
    = if null(x) then cons(i, nil)
      else let s = least(x)
        in if i < s then cons(i, sort(x))
          else Go[[cons(s, sort(if s = i then x
            else cons(i, rest(x, s)))] I3]] I3
        end
    = if null(x) then cons(i, nil)
      else let s = least(x)
        in if i < s then cons(i, sort(x))
          else Go[[if s = i then cons(s, sort(x))
            else cons(s, sort(cons(i, rest(x, s)))] I3]] I3
        end
    = if null(x) then cons(i, nil)
      else let s = least(x)
        in if i < s then cons(i, sort(x))
          else if s = i then cons(s, sort(x))
            else cons(s, sort(cons(i, rest(x, s))))
        end
  ≡ eo

```

(f_e, if₂, if₂₃)*

(f_e, if₂, let_{1-v}, let_{1-v})*

(f_e, if₂)*

(p_n, f_n, p₂, f₁)

(f_e)*

(p_e, if₂₃, p_n, f_n, v_n)*

(f_{1-if}, p_{2-if})*

(if₂₃, p_n, f_n, v_n)

Step 2. Define $sort^*(x, i, r)$, where $r = sort(x)$, to be e^o with $sort(cons(i, rest(x, s)))$ in e^o replaced by $sort^*(rest(x, s), i, sort(rest(x, s)))$. We get

```

sort*(x, i, r)
= Dt[[eo] ∅]
= if null(x) then cons(i, nil)
  else let s = least(x)
    in if i < s then cons(i, sort(x))
      else if s = i then cons(s, sort(x))
        else cons(s, sort*(rest(x, s), i, sort(rest(x, s))))
    end
≡ e*

```

Step 3. First, we itemize here those G_a transformations that are needed.

Step 4. Apply standard optimizations on $sort^*(x, i, r) = e^t$. The function definitions $sort$, $least$, $rest$, and the argument x to the function $sort^*$ are eliminated eventually.

```
(merge conditional branches)
sort*(x, i, r)
= if null(r) then cons(i, nil)
  else let s = car(r)
        in if i ≤ s then cons(i, r)
            else cons(s, sort*(rest(x, s), i, cdr(r)))
        end
(eliminate redundant argument x to sort* to get sort')
sort'(i, r)
= if null(r) then cons(i, nil)
  else let s = car(r)
        in if i ≤ s then cons(i, r)
            else cons(s, sort'(i, cdr(r)))
        end
(miscellaneous: let can be either saved or not)
sort'(i, r)
= if null(r) then cons(i, nil)
  else if i ≤ car(r) then cons(i, r)
        else cons(car(r), sort'(i, cdr(r)))
```

Analysis. For x of length n , the time complexity of $sort(x)$ is $O(n^2)$, and thus $sort(cons(i, x))$ takes $O(n^2)$ time. But $sort'$ inserts the element i into the sorted list r in $O(n)$ time. Thus, we have obtained a linear speedup by making use of the cached result r .

7 Discussion

In our presentation, we have emphasized conveying basic ideas and exposing basic problems. At several places in the transformation, algorithms need to be further studied in order to make the transformation approach efficient. Moreover, effective algorithms are also important for achieving a higher degree of incrementality in the derived programs, which is the main concern of this paper.

Issues related to the degree of incrementality. Informally, we say that a derived program has a higher degree of incrementality if more of its computation is avoided by using the cached result. There are two main issues in the transformation approach that relate to the degree of incrementality achievable in derived programs.

First, unfolding strategies affect the degree of incrementality that can be derived. Since unrestricted unfoldings may lead to nontermination, conditions are attached to unfolding rules. However, at the same time that these conditions restrict unfoldings, they may also mask opportunities for discovering a higher degree of incrementality. In \mathcal{G}_o , we only unfold calls to a function at most once for each of its favorable argument patterns. Thus, some computations depending purely on x may not be separated out in Step 1 and hence can not be replaced by their substitution candidates in Step 3. In \mathcal{G}_a , we only unfold function calls that are expandable. As a result, some auxiliary

residual expressions with substitution candidates may not be found by \mathcal{G}_a . Therefore, in both transformations, more lenient unfolding conditions could be studied to help achieve a higher degree of incrementality in the derived programs.

Second, heuristics for time comparisons affect the effectiveness of the transformations and hence the degree of incrementality that can be discovered. At several places in the transformations, especially in the \mathcal{S}_t transformation, we run into situations where we need to decide whether $t(e_1) \leq t(r_1)$ or $t(e_1) \leq t(r_1)$, and if so, replace e_1 with r_1 . A quick heuristic may simply give negative answers to most comparisons. Thus, no replacement would occur in these situations and there would be less use of cached results. So good timing heuristics are important for the derivation of a higher degree of incrementality. Note, however, that a more accurate heuristic may take a long time, making the transformation inefficient. Therefore, fast heuristics are also important for the efficiency of the transformation.

On the other hand, these issues are of concern because the problems are intractable in the general setting. If we give up generality, a lot of work can be done to frame the classes of domain problems that need to be solved and/or the features of programming languages that can be used for coding the domain problems. Then we can study efficient algorithms for these frameworks without sacrificing the degree of incrementality that can be derived. Of course, relating this potential work with existing work in incremental computation is itself another important task.

Applicability and limitation. The transformation approach presented here can only derive a set F' of function definitions such that in the set F' , f'_0 is a differentiated version of f_0 , and each function other than f_0 in the original set F is either eliminated in Step 4 (if all function calls to them are either unfolded in Step 1 or substituted out in Step 3) or kept unchanged in F' .

The reason for this limitation is as follows: Basically, \mathcal{G}_o of Step 1 naively decides whether a function application is unfolded or not using rule f_e . For any function, if rule f_e does not apply at one of its calls and this function application does not have a substitution candidate, then the definition of this function is needed as a whole and is kept unchanged henceforth. During the whole transformation, the function f_0 is derived towards its differentiated version f'_0 under the change \oplus . But this is only made possible through the introduction of the function f^* with the application of transformation \mathcal{D}_t . For any function other than f_0 , say g , if g is defined recursively and there exists some function applications of g that do not have substitution candidates, then without an application of transformation \mathcal{D}_t on g , its definition has to be kept unchanged for recursive calls to itself. In summary, the limitation of the approach above is the result of the naive decision about function unfoldings by the rule f_e of \mathcal{G}_o as well as the naive separation and sequential execution of

the four basic steps as presented above.

Overcoming the limitation. Suppose we are given a function f_0 in the context of a set F of function definitions and an input change operation \oplus , and we want to derive a differentiated version f'_0 of f_0 in the context of a set F' of function definitions. In general, any function g in F may have its differentiated versions in F' from the derivation. These differentiated versions may also be incremental versions of g under certain change operations caused by \oplus . In order to obtain such incrementality in the derived programs, we need to modify the basic unfolding rules and combine the four basic steps in a more organic way.

We propose the following idea of *focusing* and *switching*: During the transformation, the focus is on the function application that is presently transformed, say the function is g , and the goal is to apply the four-step transformation to this application of g . If there is no function call in g , then, all four steps can be applied subsequently and we are done with this application of g . Otherwise, if g has a function call at some point, say to the function h , then we first apply the four-step transformation on the part of g before the call point to h , then we switch attention to the application of h , which means that the goal now is to focus on this application of h and apply the four-step transformation on it. This procedure is applied recursively to all function calls. When we are done with this application of h , we switch our focus back from h to g to continue the transformation on g . We apply the four-step transformation on the part of g after the call point to h until we meet another function call or we finish the transformation on this application of g . Of course, we start with function f_0 .

To implement the above approach, we propose that the four-step transformation presented in Section 5 be modified as follows: Associate with each function h in F a list of its differentiated versions, initially empty. These differentiated versions are introduced at some calls to h that depend on both x and y . The transformation works as follows: At a call $h(e_1, \dots, e_n)$ that depends on both x and y , if we can not use any of h 's differentiated versions that have previously been introduced in h 's list (in a fashion similar to the way that \mathcal{D}_t uses its rule f_s to replace a call to f with a call to f^*), then we introduce a new function definition h^* for this call (in a fashion similar to the way Step 2 introduces f^*), add h^* to the list of differentiated versions of h , and proceed to transform on the unfolded function body of this h^* . Of course, more special care is needed to deal with recursive function applications.

A thorough presentation of the focusing-and-switching technique sketched here is forthcoming.

8 Related Work

A comprehensive guide to the literature on incremental computation has appeared in [23]. Despite the relatively diverse categories discussed in [23], most of the work can be divided into three classes.

The first class includes particular incremental algorithms that deal with particular input changes to particular problems. Examples are incremental attribute evaluation [26] [39] [12], incremental data-flow analysis [29] [28], incremental circuit evaluation [2], incremental shortest path problems [22], etc. The study of dynamic algorithms [40] [8] [24] can be viewed as falling into this class. Although efforts in this class are directed towards particular incremental algorithms, they apply to a broad class of problems, e.g., any attribute grammar, any circuit, etc.

In the second class, rather than manually developing particular incremental algorithms for particular problems, general evaluation frameworks are pursued for achieving incremental computation automatically, e.g., incremental computation via function caching [21], formal program manipulations like traditional partial evaluation [36] [35] [9], and incremental computation as a program abstraction [11], etc. Note that, when attribute grammars are used as a general description language for applications, the attribute evaluation framework can be viewed as falling into this class too. The characteristic of work in this class is that a general incremental evaluation model is run for any given application program and no explicitly incremental version of the program is derived and run autonomously by a standard evaluator. For this reason, these solutions to the incremental computation problem for particular applications are not readily comparable with explicitly derived incremental algorithms such as those in the first class [36] [20].

In the third class, explicitly incremental programs are derived from non-incremental programs by formal transformation techniques such as finite differencing [18] [41]. Typically, programs are written in very high-level languages with aggregate data structures, e.g., sets and bags, and fixed rules are offered for transforming aggregate operations into more efficient incremental operations. What is not provided is a method for deriving an incremental program from a non-incremental program written in standard programming languages like Pascal or Lisp.

Our work is closest in spirit to the finite differencing techniques of the third class. The name “finite differencing” was originally given by Paige and Koenig [14] [18]. Their work generalizes Cocke’s method of strength reduction [6] and provides a convenient framework with which to implement a host of program transformations including Earley’s “iterator inversion” [7]. They develop a set of rules for differentiating set-theoretic expressions and combine them using a chain rule to derive inexpensive incremental programs. Such techniques are indispensable as part of an optimizing compiler for programs written in very high-level languages like SETL or APL [17] [5]. The APTS

[16] program transformation system has been developed for such purposes. Our technique differs from that work in that it applies to programs written in a standard language like Lisp. In general, such programs are written at a lower abstraction level so that a fixed set of rules for differentiating expressions involving complex objects like sets is not sufficient. The technique we propose can be regarded as a principle and a systematic approach. Using this approach, incrementalities can be *discovered* in existing programs written in standard languages.

KIDS [32] [33] is a semiautomatic program development system that aims to derive programs from high-level specifications, as is APTS. The system performs algorithm design [34], deductive inference, program simplification, partial evaluation, finite differencing, data type refinement, etc. Its version of finite differencing is developed for the optimization of its derived functional programs and it has two basic operations: abstraction and simplification. Abstraction of a function $F(x)$ with respect to $E(x)$ results a function $F(x, c)$ with the invariance $c = E(x)$ and changes function calls $F(U(x))$ to $F(U(x), E(U(x)))$. Simplification uses the invariance to simplify F by replacing occurrences of $E(x)$ with c and, in particular, aims to compute $E(U(x))$ using the invariance $c = E(x)$. As is pointed out in [32] [33], the real benefit of this optimization comes from the potential use of $c = E(x)$ in the computation of $E(U(x))$. However, no special technique is provided except the observation that distributive laws can often be applied to $E(U(x))$ yielding $U'(E(x))$ and then $U'(c)$. Note that the essence of this optimization is to compute E (used by F) incrementally under the change U (caused by F), which is exactly what our work is aimed at. Also note that $U(x)$ is not an accurate description of the change expression, which can actually have other free variables, say y . As far as we know, KIDS has no formal notion of incremental programs of $E(x)$ with respect to $U(x, y)$, as we give with our notation for $f(x)$ with respect to $x \oplus y$, and, other than the distributive laws, offers no special treatment for achieving incrementalities of $E(U(x, y))$ using $c = E(x)$, as we do for computing $f(x \oplus y)$ using $r = f(x)$ with our transformations \mathcal{G}_o , \mathcal{D}_t , \mathcal{S}_t and \mathcal{G}_a . The Munich CIP project [19] has a finite differencing strategy similar to that of KIDS.

Each of the three classes of work discussed above emphasizes a different, although important, aspect of incremental computation. But none can be regarded as a general model that subsumes the others. A general model \mathcal{M} for incremental computation would be one that, given a non-incremental program f written in some language \mathcal{L} and an input change \oplus (which is also describable in \mathcal{L}), derives f' , an incremental version of f under \oplus . Such an \mathcal{M} can be viewed as a general model that addresses all three approaches discussed above for the following reasons: First, the development of particular incremental algorithms in the first class is a special case of \mathcal{M} , where f and \oplus are fixed according to the particular problems being studied. Secondly, work in the second class can be

regarded as a specialization of \mathcal{M} to deal with a special class of input change, namely, the change operation \oplus that always takes an old input x and returns a new input x' with part of x retained in x' and the rest of x' being new parameters y . Thirdly, work in the third class is exactly the study of \mathcal{M} where the language \mathcal{L} for describing the programs and the changes is limited to very high-level languages.

Our work attacks the problem of deriving incremental programs from non-incremental programs written in a standard programming language. It begins the study of a general model for incremental computation along lines unique and distinct from all other approaches. Although this problem is, in general, very hard, we have shown that effective approaches can be developed to derive efficient incremental programs by combining particular program transformation techniques. By studying general techniques, we aim to better understand the essence of incremental computation. We also aim to establish a general framework in which different ideas on incremental computation can be integrated. By specializing the general techniques to different applications, we will be able to obtain particular incremental algorithms, particular incremental computation techniques, and particular incremental computation languages. Their applications could include most problems discussed in the literature [23] on incremental computation.

References

- [1] A. V. Aho, R. Sethi, and J. D. Ullman. *Compilers, Principles, Techniques, and Tools*. Addison-Wesley series in Computer Science. Addison-Wesley Publishing Company, 1986.
- [2] B. Alpern, R. Hoover, B. Rosen, P. Sweeney, and K. Zadeck. Incremental evaluation of computational circuits. In *Proceedings of the First Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 32–42, San Francisco, California, January 1990.
- [3] R. A. Ballance, S. L. Graham, and M. L. V. D. Vanter. The *Pan* language-based editing system. *ACM Transactions on Programming Languages and Systems*, 1(1):95–127, January 1992.
- [4] B. Bjørner, A. P. Ershov, and N. D. Jones, editors. *Partial Evaluation and Mixed Computation*. North-Holland, 1988. Proceedings of the IFIP TC2 Workshop on Partial Evaluation and Mixed Computation, Gammel Avernæs, Denmark, October 1987.
- [5] J. Cai and R. Paige. Program derivation by fixed point computation. *Science of Computer Programming*, 11:197–261, September 1988/89.

- [6] J. Cocke and K. Kennedy. An algorithm for reduction of operator strength. *Communications of the ACM*, 20(11):850–856, November 1977.
- [7] J. Earley. High level iterators and a method for automatically designing data structure representation. *Journal of Computer Languages*, 1:321–342, 1976.
- [8] D. Eppstein, Z. Galil, G. Italiano, and A. Nissenzweig. Sparsification - a technique for speeding up dynamic graph algorithms. In *Proceedings of the 33rd Annual IEEE Symposium on FOCS*, Pittsburgh, Pennsylvania, October 1992.
- [9] J. Field and T. Teitelbaum. Incremental reduction in the lambda calculus. In *Proceedings of the ACM '90 Conference on LISP and Functional Programming*, pages 307–322, 1990.
- [10] Y. Futamura and K. Nogi. Generalized partial evaluation. In *Partial Evaluation and Mixed Computation*, pages 133–151. North-Holland, 1988. In [4].
- [11] R. Hoover. Alphonse: Incremental computation as a programming abstraction. In *Proceedings of the ACM SIGPLAN '92 Conference on PLDI*, pages 261–272, California, June 1992.
- [12] P. Lipps, U. Möncke, M. Olk, and R. Wilhelm. Attribute (re)evaluation in OPTRAN. *Acta Informatica*, 26:213–239, 1988.
- [13] Y. A. Liu. Theoretical results on deriving incremental programs. Unpublished notes, 1992.
- [14] R. Paige. *Formal Differentiation: A Program Synthesis Technique*, volume 7 of *Computer Science. Artificial Intelligence*. UMI Research Press, Ann Arbor, Michigan, 1981. Revision of Ph.D. Thesis - New York University, 1979.
- [15] R. Paige. Transformational programming – applications to algorithms and systems. In *Conference Record of the 10th Annual ACM Symposium on POPL*, pages 73–87, January 1983.
- [16] R. Paige. Symbolic finite differencing - part I. In *Proceedings of the 3rd ESOP*, pages 36–56, Copenhagen, Denmark, May 1990. Springer-Verlag. LNCS 432.
- [17] R. Paige and F. Henglein. Mechanical translation of set theoretic problem specifications into efficient ram code - a case study. *Journal of Symbolic Computation*, 4(2):207–232, 1987.
- [18] R. Paige and S. Koenig. Finite differencing of computable expressions. *ACM Transactions on Programming Languages and Systems*, 4(3):402–454, July 1982.

- [19] H. A. Partsch. *Specification and Transformation of Programs - A Formal Approach to Software Development*. Texts and Monographs in Computer Science. Springer-Verlag, 1990.
- [20] W. Pugh. Comments on ‘incremental computation via partial evaluation’. Private communications, March 1991.
- [21] W. Pugh and T. Teitelbaum. Incremental computation via function caching. In *Conference Record of the 16th Annual ACM Symposium on POPL*, pages 315–328, January 1989.
- [22] G. Ramalingam and T. Reps. On the computational complexity of incremental algorithms. Technical Report TR-1033, Computer Sciences Department, University of Wisconsin, Madison, Wisconsin, August 1991.
- [23] G. Ramalingam and T. Reps. A categorized bibliography on incremental computation. In *Conference Record of the 20th Annual ACM Symposium on POPL*, pages 502–510, Charleston, South Carolina, January 1993.
- [24] M. Rauch. Fully dynamic biconnectivity in graphs. In *Proceedings of the 33rd Annual IEEE Symposium on FOCS*, pages 50–59, Pittsburgh, Pennsylvania, October 1992.
- [25] T. Reps and T. Teitelbaum. *The Synthesizer Generator: A System for Constructing Language-Based Editors*. Springer-Verlag, New York, 1988.
- [26] T. Reps, T. Teitelbaum, and A. Demers. Incremental context-dependent analysis for language-based editors. *ACM Transactions on Programming Languages and Systems*, 5(3):449–477, July 1983.
- [27] M. Rosendahl. Automatic complexity analysis. In *Proceedings of the 4th International Conference on FPCA*, pages 144–156, London, September 1989.
- [28] B. Ryder, T. Marlowe, and M. Paull. Conditions for incremental iteration: examples and counterexamples. *Science of Computer Programming*, 11(1):1–15, 1988.
- [29] B. Ryder and M. Paull. Incremental data flow analysis algorithms. *ACM Transactions on Programming Languages and Systems*, 10(1):1–50, January 1988.
- [30] D. Sands. Complexity analysis for a lazy higher-order language. In *Proceedings of the 3rd ESOP*, pages 361–376, Copenhagen, Denmark, May 1990. Springer-Verlag. LNCS 432.
- [31] M. Sharir. Some observations concerning formal differentiation of set theoretic expressions. *ACM Transactions on Programming Languages and Systems*, 4(2):196–225, April 1982.

- [32] D. R. Smith. KIDS: A semiautomatic program development system. *IEEE Transactions on Software Engineering*, 16(9):1024–1043, September 1990.
- [33] D. R. Smith. KIDS - a knowledge-based software development system. In M. R. Lowry and R. D. MacCartney, editors, *Automating Software Design*, chapter 19, pages 483–514. AAAI Press/The MIT Press, 1991. Proceedings of the Workshop on Automating Software Design, AAAI '88.
- [34] D. R. Smith and M. R. Lowry. Algorithm theories and design tactics. *Science of Computer Programming*, 14:305–321, 1990.
- [35] R. Sundaresh. Building incremental programs using partial evaluation. In *Proceedings of the ACM SIGPLAN Symposium on PEPM*, pages 83–93, Yale University, June 1991.
- [36] R. Sundaresh and P. Hudak. Incremental computation via partial evaluation. In *Conference Record of the 18th Annual ACM Symposium on POPL*, pages 1–13, January 1991.
- [37] A. Takano. Generalized partial computation for a lazy functional language. In *Proceedings of the Symposium on PEPM*, pages 1–11, Yale University, June 1991.
- [38] B. Wegbreit. Mechanical program analysis. *Communications of the ACM*, 18(9):528–538, September 1975.
- [39] D. Yeh and U. Kastens. Improvements on an incremental evaluation algorithm for ordered attribute grammars. *SIGPLAN Notices*, 23(12):45–50, 1988.
- [40] D. M. Yellin. Speeding up dynamic transitive closure for bounded degree graphs. *Acta Informatica*, 30(4):369–384, July 1993.
- [41] D. M. Yellin and R. E. Strom. INC: A language for incremental computations. *ACM Transactions on Programming Languages and Systems*, 13(2):211–236, April 1991.