

ARITY RAISER AND ITS USE IN PROGRAM SPECIALIZATION

Sergei A. Romanenko

Keldysh Institute of Applied Mathematics
Academy of Sciences of the USSR
Miusskaya Sq.4, SU-125047, Moscow, USSR

Experiments on generating compilers by specializing specializers with respect to interpreters have shown that the compilers thus obtained have a natural structure only if the specializer does *variable splitting*. Variable splitting can result in a residual program using several variables to represent the values of a single variable of the original program. In the case of functional programming variable splitting is done by raising the arities of functions. The paper describes the structure and principles of operation of an arity raiser dealing with programs in a subset of pure Lisp.

Keywords: arity raiser, compiler generator, partial evaluation, retyping, specializer, variable splitting.

INTRODUCTION

Program specialization [Dixon 71] seems to be a promising and powerful technique that can lead to new program development methodology.

By *program specialization* we understand constructing, when given a "general-purpose" program and some restriction on its usage, a more efficient "specialized" residual program. Being optimized and simplified version of the original program, the residual program, however, must be equivalent to the original one when used according to the restriction. By *specializer* we understand a system that, given a program and a restriction, will produce a specialized version of the original program.

Program specialization can be achieved by making use of different techniques, such as *driving* [Turchin 72], *fold-unfold method* [Burstall 77], *partial evaluation* [Futamura 71], [Beckman 76], *mixed computation* [Ershov 78], [Bulyonkov 84], the *analysis of computational configurations* [Turchin 79], [Turchin 86], *variable splitting* [Sestoft 86], and *arity raising* [Romanenko 88].

The above techniques deal, for the most part, with two problems: *control restructuring* and *data retyping* (i.e. changing representation of data).

As far as the control restructuring is concerned, various specialization techniques differ in the extent to which the program is reorganized.

In the case of *monovariant* specialization any control point in the original program gives rise to zero or one control point in the residual program.

In the case of *polyvariant* specialization a control point can give rise to more than one control point in the residual program.

In the case of *monogenetic* specialization any control point in the residual program is produced from a single control point of the original program.

In the case of *polygenetic* specialization a control point in the residual program may be produced from several control points of the original program.

As far as the data representation is concerned, various specialization techniques differ in the use they make of *retyping*.

Driving [Turchin 72] and the analysis of configurations [Turchin 79], [Turchin 86], which deal with functional programs, can be classified as polyvariant polygenetic methods with retyping.

Monovariant monogenetic techniques for imperative programs are studied in [Ershov 78]. Papers [Bulyonkov 84], [Barzdin 88] concern polyvariant monogenetic specialization techniques for imperative programs.

The *transformational approach* [Ershov 81], [Ostrovski 88] is believed to include, at least potentially, all conceivable techniques of program specialization, not excluding the polygenetic ones.

Of course, the more powerful techniques tend to be rather expensive, and it is difficult to make them completely automatic. Thus the choice of appropriate specialization techniques depends on the class of problems to be solved.

An interesting application of specializers is compiler generation. It was found by Y. Futamura [Futamura 71] that interpreters can be converted to compilers by specializing a specializer with respect to the interpreters. Several years later it was realized [Beckman 76] that a transformer of interpreters into compilers can be produced by specializing a specializer with respect to a specializer.

To put this approach into practice, we have to overcome the following difficulty. On the one hand, the specializer has to be sophisticated enough

to achieve non-trivial specialization. On the other hand, to be specializable, the specializer can't afford to be too complicated.

The group under N.D.Jones at Copenhagen university was the first to overcome the above difficulty [Jones 85], [Sestoft 86], [Sestoft 88]. Since experiments had shown the monovariant specialization to be unsatisfactory for this application, the specializer had to do the polyvariant specialization. Again, the monogenetic specialization proved to be adequate for the purpose (despite there being a lot of problems that have to be dealt with by polygenetic specialization [Turchin 82], [Wadler 88]).

The usefulness of retying proved to be more problematic. It was found that retying can be dispensed with at the cost of the residual programs having rather unnatural structure. Suppose, for example, that an interpreter is to be specialized with respect to a program. Since the interpreter is supposed to accept an arbitrary input program, the number of variables in this program cannot be known in advance. Thus the variable's values are likely to be represented in the interpreter as a single value assigned to one of the interpreter's variables. If the specializer is unable to split this variable, the residual program will use a single variable to represent all the values. A reasonable residual program, however, would keep each value in a separate variable [Sestoft 86].

To rectify the drawback, the author suggested that the Copenhagen specializer should be supplemented with an additional phase, whose purpose would be to do variable splitting [Romanenko 88]. In the case of a functional language, variable splitting reduces to increasing the number of functions' parameters, for which reason this additional phase was given the name *arity raiser*. As pointed out by T.Mogensen arity raising is just a special case of retying, thus any arity raiser is a retyper.

The arity raiser was found to improve the structure of residual programs without making the specializer excessively slow and intricate.

The alternative to the arity raiser is to split variables *on-line*, i.e. at the time the residual program is being generated [Turchin 86], [Mogensen 88]. This approach, however, can result in a mammoth, sluggish specializer.

A short description of the ideas behind the arity raiser can be found in [Romanenko 88]. The present paper gives a detailed account of the structure and principles of operation of an arity raiser dealing with programs in a subset of pure Lisp.

1. THE LANGUAGE MIXWELL

In the following we consider programs written in the language Mixwell, which is a small subset of pure Lisp and was used as the subject language in the Copenhagen specializer MIX [Sestoft 86]. Here is Mixwell's abstract syntax.

pgm	∈ Program	programs
fd	∈ FnDef	function definitions
exp, e	∈ Exp	expressions
f	∈ FName	function names
x	∈ VName	variable names
\mathcal{A}	∈ Atom	Lisp atoms
\mathcal{E}	∈ SExp	Lisp S-expressions

```

pgm ::= fd1; ... fdn;
fd   ::= f(x1, ..., xm) = exp
exp ::= x | quote  $\mathcal{E}$  | if exp0 then exp1 else exp2 | call f(exp1, ..., expm)
      | car(exp) | cdr(exp) | cons(exp1, exp2) | atom(exp) | equal(exp1, exp2)
 $\mathcal{E}$  ::=  $\mathcal{A}$  | ( $\mathcal{E}_1$  .  $\mathcal{E}_2$ )

```

A Mixwell program is a list of function definitions, the first function being the goal function. The goal function is to be called first, and inputs to the program are through the parameters of this function.

The body of a function is an expression, which is constructed from variables appearing in the function's formal parameter list, from constants `quote` and operators `car`, `cdr`, `cons`, `atom` and `equal` (as in Lisp), conditionals `if` and defined function calls `call`.

The only data type is well-founded (i.e. non-circular) S-expressions as known from Lisp.

All primitive and defined functions, except the conditional `if`, are strict in all positions. All parameters are called by value.

We use some "syntactic sugar". The keyword `call` is omitted in cases where the name of the function being called is different from the names of the primitive functions. `quote \mathcal{E}` can be written as `' \mathcal{E}` , `cons(exp1, exp2)` as `exp1 :: exp2`, `equal(exp1, exp2)` as `exp1 = exp2`. Constants (`\mathcal{E}_1 . (\mathcal{E}_2 (\mathcal{E}_n . nil) ...)`) can be written as `(\mathcal{E}_1 \mathcal{E}_2 ... \mathcal{E}_n)`.

2. SPLITTING A FORMAL PARAMETER

Suppose the definition of function `f` in a program has the form

$f(\dots, x_k, \dots) = \text{exp}$. Then the following transformation will be referred to as *the splitting of the function's k-th parameter*.

Let x' and x'' be two variables different from all formal parameters of the function f . Then the splitting of x into x' and x'' can be done in two steps.

At the first step, the original definition of f is replaced with

$$f(\dots, x', x'', \dots) = \text{exp}[x_k \rightarrow \text{cons}(x', x'')]$$

where $\text{exp}[x_k \rightarrow \text{cons}(x', x'')]$ denotes the expression obtained from exp by replacing all occurrences of x_k with $\text{cons}(x', x'')$.

At the second step, all calls of the function f in all function definitions are transformed, each call of the form $\text{call } f(\dots, e_k, \dots)$ being replaced with $\text{call } f(\dots, \text{car}(e_k), \text{cdr}(e_k), \dots)$.

Thus, the original variable x_k is replaced by two new variables x' and x'' containing enough information for the value of x_k , if needed, to be reconstructed. To put it more exactly, the value of x_k can be obtained by evaluating the expression $\text{cons}(x', x'')$.

The fact that the formal parameter x of the function f is to be split into two variables x' and x'' will, for the brevity's sake, be written as $f(x \rightarrow x' :: x'')$.

Example. Consider the program $f(x) = g(x :: x); g(u) = \text{cdr}(u);$. By $g(u \rightarrow u1 :: u2)$ we get the program $f(x) = g(x :: x); g(u1, u2) = \text{cdr}(u1 :: u2);$. Then, by splitting the argument in the calls of g we get $f(x) = g(\text{car}(x :: x), \text{cdr}(x :: x)); g(u1, u2) = \text{cdr}(u1 :: u2);$.

This program can be locally optimized, which results in $f(x) = g(x, x); g(u1, u2) = u2;$. Now we see that variable splitting is capable of producing parameters whose values are certain not to be needed. Such parameters can be recognized by a kind of backward analysis [Hughes 88] and eliminated. In the above program we can remove the parameter $u1$ of the function g , which gives the program $f(x) = g(x); g(u2) = u2;$.

Thus, the principal use of variable splitting consists in paving the way for other transformations such as local optimization and elimination of unneeded parameters, the latter being, in a sense, a kind of "garbage collection at compile time".

3.CONDITIONS OF THE VARIABLE SPLITTING CORRECTNESS

The program transformation described above can be incorrect. For

example, after $g(u \rightarrow u1 :: u2)$ for the program $f(x) = g('a); g(u) = u;$ we get $f(x) = g(car('a),cdr('a)); g(u1,u2) = u1 :: u2;.$

It is evident that the transformed program is not equivalent to the original one, because the original program terminates, with the result being the atom 'a, whereas the transformed program fails to apply *car* or *cdr* to the atom 'a and terminates abnormally. Thus we come to the conclusion:

Before splitting a parameter, we must make sure that, when the program is run, it is impossible for the parameter's value to be an atom!

Hence, to split a variable, we need to have a description of the structure of its values. Such descriptions will be referred to as *types* of variables.

4. ANALYSIS OF RUN TIME TYPES

To describe the structure of values to be taken by a variable, we use the following set of types.

t	\in Type	types
A	\in Atom	Lisp atoms

$t ::= any \mid atom(A) \mid cons(t_1, t_2) \mid \perp$

We assume the set of types to be equipped with reflexive partial ordering \leq recursively defined by the following rules:

- (i) $\perp \leq t \leq any$ for all types t .
- (ii) $cons(t'_1, t'_2) \leq cons(t''_1, t''_2)$ if $t'_1 \leq t''_1$ and $t'_2 \leq t''_2$.

If $t' \leq t''$ and $t' \neq t''$, the type t'' is said to be *more general* than the type t' .

The set of types is a lattice, as for all types $t', t'' \in Type$ there exist their least upper bound $t' \sqcup t''$ and their greatest lower bound $t' \sqcap t''$. Each set of types $TeP(Type)$ has its least upper bound $\sqcup T$. Thus the set of types is a pointed continuous partial ordering (CPO) with the bottom \perp [Schmidt 86]. It can be easily seen that the set of types has no chains of infinite height. In addition, each *finite* $TeP(Type)$ has its greatest lower bound $\sqcap T$.

A type represents a set of S-expressions. More specifically, let us define an "abstraction" function Abs mapping sets of S-expressions into types. Abs is defined in terms of an auxiliary function Abs' mapping S-expressions into types.

$$\begin{aligned} Abs &\in \mathcal{P}(SExp) \rightarrow Type \\ Abs' &\in SExp \rightarrow Type \\ Abs[E] &= \sqcup \{Abs'[\mathcal{E}] \mid \mathcal{E} \in E\} \\ Abs'[A] &= atom(A) \\ Abs'[(\mathcal{E}' . \mathcal{E}'')] &= cons(Abs'[\mathcal{E}'], Abs'[\mathcal{E}'']) \end{aligned}$$

Let us define a "concretization" function Co reconstructing the set of S-expressions from a type:

$$\begin{aligned} Co &\in Type \rightarrow \mathcal{P}(SExp) \\ Co[any] &= SExp \\ Co[atom(A)] &= \{A\} \\ Co[cons(t', t'')] &= \{(\mathcal{E}' . \mathcal{E}'') \mid \mathcal{E}' \in Co[t'] \text{ and } \mathcal{E}'' \in Co[t'']\} \\ Co[\perp] &= \{\} \end{aligned}$$

The following relations hold: $Abs[Co[t]] = t$ and $E \subseteq Co[Abs[E]]$.

Now let x be a variable in a program. The problem is to find a type t such that $\mathcal{E} \in Co[t]$ for all \mathcal{E} that can be taken as value by x when the program is run. It can be done by *abstract interpretation* [Jones 86] of the program, which amounts to performing the program's computations using abstract values in place of the actual ones.

[Omitted: the type analysis algorithm.]

The type analysis above can, in a sense, be regarded as a monovariant, monogenetic version of the "configuration analysis" as used in the Supercompiler [Turchin 89], [Turchin 86].

5. USING TYPE INFORMATION FOR VARIABLE SPLITTING

The variable splitting transformation as described above splits only one of a function's parameters. However, the information provided by an argument type description is sufficient for all function's parameters to be split at once.

Suppose the type t assigned to a variable x contains some occurrences of the type any , which will be referred to as "gaps".

It is obvious that all values of the variable x can be different only at places corresponding to the gaps, and must be congruent at all other

places. Therefore, if the type t contains m gaps, any S-expression $\mathcal{E} \in \text{Co}[t]$ is completely determined by its parts corresponding to the gaps in the type t . This enables the variable x to be *retyped* by replacing it with m new variables, which are to be assigned the parts of the variable's values corresponding to the gaps.

For example, if a function's parameter x has the type $\text{cons}(\text{cons}(\text{any}, \text{atom}(a)), \text{any})$, then x can be represented by two new parameters u_1 and u_2 , in which case all occurrences of x in the function's body must be replaced with the expression $(u_1 :: 'a') :: u_2$.

[Omitted: the algorithm for splitting parameters according their types.]

6. CODE DUPLICATION RISK

Example. Consider the program $f(z) = \text{swap}(\text{unzip}(z, 'nil, 'nil)); \text{swap}(v) = \text{cdr}(v) :: \text{car}(v); \text{unzip}(u, x, y) = \text{if } u = 'nil \text{ then } x :: y \text{ else } \text{unzip}(\text{cdr}(u), \text{car}(\text{car}(u)) :: x, \text{cdr}(\text{car}(u)) :: y);$.

Any result produced by the function `unzip` is of the type $\text{cons}(\text{any}, \text{any})$, hence this type can be assigned to the parameter v of the function `swap`. Thus we can perform $\text{swap}(v \rightarrow v_1 :: v_2)$. But this gives rise to two copies of the expression $\text{unzip}(z, 'nil, 'nil)$, which is bad for two reasons. First, duplicating expressions can result in huge programs being produced. Second, code duplication can lead to repeated evaluation of expressions. Both of the problems arise in the above example.

The risk of code duplication and repeated evaluation can be avoided by following the principle of *selector non-introduction*:

All selectors produced by variable splitting must be eliminable by means of local optimization.

What is the drawback of the type analysis described above? The point is that this analysis tells us whether a selector in the program is certain to be applicable at run time, whereas we need to know whether the selector can be applied *symbolically* at the time the program is being optimized.

The feasibility of the symbolic application of a selector to the expression exp , obviously, depends upon the structure of the expression itself, rather than on the structure of the result to be produced by exp at run time.

If exp has the form $\text{quote } (\mathcal{E}' . \mathcal{E}'')$, the symbolic application is feasible, $\text{car}(\text{exp})$ being reducible to $\text{quote } \mathcal{E}'$, and $\text{cdr}(\text{exp})$ being reducible to $\text{quote } \mathcal{E}''$.

If exp has the form $\text{exp}' :: \text{exp}''$, the symbolic application is feasible, $\text{car}(\text{exp})$ being reducible to exp' , and $\text{cdr}(\text{exp})$ being reducible to exp'' .

On the other hand, if exp has the form $\text{if } \text{exp}_0 \text{ then } \text{exp}' \text{ else } \text{exp}''$ or $\text{call } f(\text{exp}_1, \dots, \text{exp}_m)$, it is impossible to make the symbolic application without code duplication.

If exp is a variable x , the symbolic application may seem to be unfeasible. Nevertheless, splitting the parameters throughout the program may result in the variable x being split into a new expression, which may enable the symbolic application.

Example. Consider the program $f(x) = g(x :: x)$; $g(u) = h(u)$; $h(v) = \text{cdr}(v)$; . After $g(u \rightarrow u_1 :: u_2)$, we get $f(x) = g(x, x)$; $g(u_1, u_2) = h(u_1 :: u_2)$; $h(v) = \text{cdr}(v)$; . Now, after $h(v \rightarrow v_1 :: v_2)$, we get $f(x) = g(x, x)$; $g(u_1, u_2) = h(u_1, u_2)$; $h(v_1, v_2) = v_2$; .

7. ANALYSIS OF OPTIMIZATION TIME TYPES

As can be seen from the above, we need to know the structure of *symbolic* values assigned to variables at the time the program is being optimized, rather than the structure of ordinary values assigned to variables at the time the program is run. Thus, what we are really interested in are the optimization time types, rather than the run time types.

To find them, we can use the same set of types as has been used for analyzing the run time types.

Suppose we have a program defining functions f_1, \dots, f_h . Let $F = \{f_1, \dots, f_h\}$, and, for each $f \in F$, $x_{f,j}$ be its j -th parameter, $a(f)$ be its arity, and body_f be its body, so that the definition of f has the form:

$$f(x_{f,1}, \dots, x_{f,a(f)}) = \text{body}_f$$

Let $\theta \in \text{Env} = \text{VName} \rightarrow \text{Type}$ be an *environment* assigning a type to each parameter of a function. Let $\alpha \in \text{ArgDescr} = F \rightarrow \text{Env}$ be an *argument type description* assigning types to each function's parameters. Let $\rho \in \text{ResDescr} = F \rightarrow \text{Type}$ be a *result type description* assigning a type to each function's result.

All the sets above are equipped with reflexive partial orderings as follows:

Env: $\theta' \leq \theta'' \Leftrightarrow \forall x \in VName \theta'(x) \leq \theta''(x)$
 ArgDescr: $\alpha' \leq \alpha'' \Leftrightarrow \forall f \in F \alpha'(f) \leq \alpha''(f)$
 ResDescr: $\rho' \leq \rho'' \Leftrightarrow \forall f \in F \rho'(f) \leq \rho''(f)$

We define two functions R and A to do the abstract interpretation using these ordered sets.

The function R, given an expression exp and an environment θ , computes the type of an expression's result.

$R \in Exp \rightarrow Env \rightarrow Type$

$R[x] \theta = \theta(x)$

$R[quote \mathcal{E}] \theta = Abs'[\mathcal{E}]$

$R[if \ exp \ then \ exp' \ else \ exp''] \theta = any$

$R[call \ f(exp_1, \dots, exp_m)] \theta = any$

$$R[car(exp)] \theta = \begin{cases} any & \text{if } R[exp] \theta = any, \\ t' & \text{if } R[exp] \theta = cons(t', t''), \\ \perp & \text{otherwise.} \end{cases}$$

$$R[cdr(exp)] \theta = \begin{cases} any & \text{if } R[exp] \theta = any, \\ t'' & \text{if } R[exp] \theta = cons(t', t''), \\ \perp & \text{otherwise.} \end{cases}$$

$R[cons(exp', exp'')] \theta = cons(R[exp'] \theta, R[exp''] \theta)$

$R[atom(exp)] \theta = any$

$R[equal(exp', exp'')] \theta = any$

The function A, given an expression exp, an environment θ , and an argument type description α , computes a new approximation to the final description of each function's parameter types.

$A \in Exp \rightarrow Env \rightarrow ArgDescr \rightarrow ArgDescr$

$A[x] \theta \alpha = \alpha$

$A[quote \mathcal{E}] \theta \alpha = \alpha$

$A[if \ exp \ then \ exp' \ else \ exp''] \theta \alpha$

$= A[exp] \theta \alpha \sqcup A[exp'] \theta \alpha \sqcup A[exp''] \theta \alpha$

$A[call \ f(exp_1, \dots, exp_m)] \theta \alpha = \alpha_{new} [f \mapsto \alpha_{new}(f) \sqcup \theta_{new}],$

where $\alpha_{new} = \sqcup \{A[exp_j] \theta \alpha\}_{j=1, \dots, m}$ and

$\theta_{new} = [x_{f, j} \mapsto R[exp_j] \theta]_{j=1, \dots, m}$

$A[car(exp)] \theta \alpha = A[exp] \theta \alpha$

$A[cdr(exp)] \theta \alpha = A[exp] \theta \alpha$

$A[cons(exp', exp'')] \theta \alpha = A[exp'] \theta \alpha \sqcup A[exp''] \theta \alpha$

$A[atom(exp)] \theta \alpha = A[exp] \theta \alpha$

$A[equal(exp', exp'')] \theta \alpha = A[exp'] \theta \alpha \sqcup A[exp''] \theta \alpha$

We want a final argument type description α that is consistent and as low as possible. This must be the least fixed point for the following

system of simultaneous equations and relations:

$$\alpha = \sqcup \{A[\text{body}_f] \alpha(f) \mid f \in F\}, \quad \alpha \geq \alpha_0$$

where α_0 is defined as follows

$$\alpha_0 = [f_1 \mapsto [x_{f_1, j} \mapsto \text{any}]_{j=1, \dots, a(f_1)}] \sqcup [f \mapsto [x_{f, j} \mapsto \perp]_{j=1, \dots, a(f)}]_{f \in F}$$

The description α_0 assigns the type *any* to the parameters of the goal function f_1 , to prevent these parameters from being split. All other parameters, on the contrary, are assigned the type \perp , there being no *a priori* information about their structure.

The least fixed point for the system above does exist because for any given program the ordered sets involved have no chains of infinite height, and the functions A and R are monotonic.

8. USEFULNESS OF VARIABLE SPLITTING

The fact that the parameters of a function f have been assigned the types t_1, \dots, t_m , for brevity's sake, will be written as $f(t_1, \dots, t_m)$.

Example. Consider the program

```
f(x) = rev(x, 'a :: 'nil);
rev(u,v) = if u = 'nil then v else rev(cdr(u), car(u) :: v);
```

The analysis of types tells us that $f(\text{any}), \text{rev}(\text{any}, \text{cons}(\text{any}, \text{any}))$. After $\text{rev}(v \rightarrow v1 :: v2)$, we get the program

```
f(x) = rev(x, 'a, 'nil);
rev(u,v1,v2) = if u = 'nil then v1 :: v2 else
               rev(cdr(u), car(u), v1 :: v2);
```

We see that the program obtained is by no means superior to the original one, because no selector has been eliminated owing to variable splitting.

Thus we see that the parameter splitting based exclusively on the information obtained by examining the structure of argument expressions, may well result in the "arity overraising", i.e. increasing the number of parameters without reducing the number of selectors in the program. The types as produced by the above analysis, describing as they do the feasibility of splitting parameters, however, provide no information on the usefulness of this splitting. The arity overraising, nevertheless, can be avoided by "adjusting" the above types in the following way.

Suppose, for example, the type t has been assigned to a parameter x . Then the splitting of the parameter can be restricted by replacing some parts of t having the form $\text{cons}(t_1, t_2)$ with any . This results in the type t being *generalized*, i.e. changed to some other type t' such that $t \leq t'$, the depth of splitting being the less the greater the type t' . Thus, for instance, the splitting $x \rightarrow x1 :: (x2 :: x3)$ corresponds to the type $\text{cons}(\text{any}, \text{cons}(\text{any}, \text{any}))$, the splitting $x \rightarrow x1 :: x2$ to the type $\text{cons}(\text{any}, \text{any})$, and no splitting to the type any .

Thus we are facing the *type generalization problem*: given a *cons* in a type, we have to decide whether this *cons* should be retained or generalized. This decision will be made on the basis of the following *selector elimination principle*:

A *cons* should be retained only if this causes a selector in the program to disappear.

Being formalized as it is, the selector elimination principle gives only an approximate description of the intuitive ideas the humans have about what does it means for a program to have a beautiful and natural structure. Nevertheless, experience has shown this principle to be likely to produce reasonable results, without any danger of the program being spoiled.

9. BACKWARD ANALYSIS

Let us consider the function definition $f(\dots, x_k, \dots) = \text{exp}$.

The k -th parameter of the function may appear at different places in the function's body exp . Is it any use splitting x_k ? To answer this question, we have to consider all occurrences of x_k in exp and to take into account their *contexts* in exp . To take an example, if exp contains the subexpression $\text{cdr}(x_k)$, it makes sense to perform the splitting $x_k \rightarrow x' :: x''$, since this will cause $\text{cdr}(x_k)$ to be replaced with $\text{cdr}(x' :: x'')$, the latter being reducible to x'' .

Example. $f(x) = g(x :: x)$; $g(u) = u$.

In this case the selector elimination principle tells us that it is no use performing the splitting $g(u \rightarrow u1 :: u2)$.

Example. $f(x) = g(x :: x); g(u) = \text{cdr}(u);$.

In this case the selector elimination principle tells us that the splitting $g(u \rightarrow u1 :: u2)$ is worth performing, since it will cause the selector cdr to disappear. And, in fact, after the splitting we get the program $f(x) = g(x, x); g(u1, u2) = u2;$.

Thus we see that the natural way of getting information about the usefulness of splitting is to make use of some kind of backward analysis [Hughes 88].

10. ACCESS PATHS AND CONTEXTS

Let exp be an expression appearing in a larger expression. We want to consider all attempts by the surrounding expression at accessing the components of exp . For example, if exp is a part of the expression $\text{car}(\text{cdr}(\text{cdr}(\text{exp})))$, then there is an attempt at accessing exp by applying selectors in the following order: cdr , cdr , car . The component to be accessed can be unambiguously identified by a sequence of selectors. This justifies the following definition.

Definition. An *access path* is a finite list (which may be empty) of selector names car and cdr .

We use the following notation. A finite list of elements a_1, \dots, a_m is written as $[a_1, \dots, a_m]$, an empty list as $[\]$. The concatenation of two lists $A = [a_1, \dots, a_m]$ and $B = [b_1, \dots, b_n]$ equal to $[a_1, \dots, a_m, b_1, \dots, b_n]$ is denoted by $A \hat{\ } B$.

The set of all access paths will be denoted by Path . Thus $\text{Path} = \{\text{car}, \text{cdr}\}^*$.

In some cases the surrounding expression tries to access several components of the expression under consideration. For this reason we have to describe the context by a set of paths, rather than by a single path.

Definition. A set of access paths $\Pi \in \mathcal{P}(\text{Path})$ is an *access context*, if it satisfies the following requirements.

- (i) $[\] \in \Pi$
- (ii) If $\pi \hat{\ } [\text{car}] \in \Pi$ or $\pi \hat{\ } [\text{cdr}] \in \Pi$, then $\pi \in \Pi$.

(i) means that an attempt at accessing the expression as a whole must be included into the context. This requirement is useful for technical

reasons. (ii) formalizes the obvious fact that a subcomponent can be accessed only by accessing the components in which the subcomponent is included.

The set of all contexts is denoted by Context.

Now consider the function definition $f(\dots, x_k, \dots) = \text{exp}$. Suppose that exp contains m occurrences of the parameter x_k in the contexts $\Pi_1, \Pi_2, \dots, \Pi_m$. What should be the total context for all occurrences of x_k ? It is clear that finding all attempts at accessing the parameter x_k amounts to finding all attempts at accessing its occurrences, thus $\Pi_1 \cup \Pi_2 \cup \dots \cup \Pi_m$ should be considered to be the total context of the parameter x_k .

11. USING CONTEXTS FOR TYPE GENERALIZATION

Let a parameter have the type t and the context Π . Then the function GenType can be easily defined which generalizes t in accordance with Π by replacing all $\text{cons}(t_1, t_2)$ unaccessed by Π with *any*.

GenType : Type \rightarrow Context \rightarrow Type

GenType[t] $\Pi = \sqcap \{ \text{GenType}'[t]\pi \mid \pi \in \Pi \}$

GenType' : Type \rightarrow Path \rightarrow Type

GenType'[any] $\pi = \text{any}$

GenType'[atom(\mathcal{A})] $\pi = \text{atom}(\mathcal{A})$

GenType'[cons(t', t'')][()] = *any*

GenType'[cons(t', t'')][[car] $\wedge \pi$] = cons(GenType'[t'] π , *any*)

GenType'[cons(t', t'')][[cdr] $\wedge \pi$] = cons(*any*, GenType'[t''] π)

GenType'[1] $\pi = 1$

It should be noted that for all $t \in \text{Type}$ and $\pi \in \Pi$ the relation $t \leq \text{GenType}'[t]\pi$ holds, therefore the set $\{ \text{GenType}'[t]\pi \mid \pi \in \Pi \}$ is finite, in spite of the fact that Π may well be infinite. Consequently, the greatest lower bound of this set does exist.

12. LATENT SELECTORS

The above considerations might have produced the expression that the context of a parameter can be determined by examining only the definition of the function concerned, without the program being globally analyzed. This is not really the case, however.

Example. $f(x) = g(x :: 'a)$; $g(u) = h(u)$; $h(v) = \text{cdr}(v)$;

The type analysis tells us that $f(\text{any})$, $g(\text{cons}(\text{any}, \text{atom}(a)))$, $h(\text{cons}(\text{any}, \text{atom}(a)))$. The variable v has the context $\{ [], [\text{cdr}] \}$. But what is the context of the variable u ? At the first glance, it may appear to be

{[]}, because there seems to be no selectors in the program attempting at accessing the variable u . Thus we, erroneously, come to the conclusion that the types should be generalized as follows: $f(\text{any})$, $g(\text{any})$, $h(\text{cons}(\text{any}, \text{atom}(a)))$. The only acceptable splitting is therefore $h(v \rightarrow v1 :: v2)$. By performing it we get $f(x) = g(x :: 'a)$; $g(u) = h(\text{car}(u), \text{cdr}(u))$; $h(v1, v2) = v2$;

This result is far from being satisfactory, because there have appeared two new selectors car and cdr , not present in the original program. This makes us draw the conclusion that the parameter access analysis has to take into account not only the selectors explicitly appearing in the program, but also the *latent selectors* to be introduced by the splitting of parameters.

Thus, if e_k is an argument expression in the function call $f(\dots, e_k, \dots)$, it would be incorrect to take its context to be $\{\{\}\}$, because there should be taken into account all attempts at accessing e_k due to the splitting of e_k . This can be done in the following way.

Let the k -th formal parameter of the function f be assigned the type t , and the total context of all its occurrences be Π . Let $t' = \text{GenType}[t]\Pi$. Then the generalized type t' gives all information about the way in which e_k is to be split. The function TypeToContext can be easily defined which converts t' into the context providing the information about all the attempts at accessing e_k due to the splitting of e_k in accordance with t' .

$\text{TypeToContext} : \text{Type} \rightarrow \text{Context}$

$\text{TypeToContext}[\text{any}] = \{\{\}\}$

$\text{TypeToContext}[\text{atom}(A)] = \{\{\}\}$

$\text{TypeToContext}[\text{cons}(t', t'')] = \{\{\}\} \cup$

$\text{car}^*(\text{TypeToContext}[t']) \cup \text{cdr}^*(\text{TypeToContext}[t''])$

$\text{TypeToContext}[_] = \{\{\}\}$

where we use the notation $\text{car}^*\Pi = \{[\text{car}]^n \mid n \in \Pi\}$, $\text{cdr}^*\Pi = \{[\text{cdr}]^n \mid n \in \Pi\}$.

Now we can determine the context of the expression e_k , assuming the k -th parameter to be assigned the type t , and the total context of all its occurrences to be Π . This context is equal to $\text{TypeToContext}[\text{GenType}[t]\Pi]$.

13. SYSTEM OF EQUATIONS FOR FINDING CONTEXTS

For each function f with the definition $f(x_{f,1}, \dots, x_{f,m}) = \text{body}_f$ let $t_{f,1}, \dots, t_{f,m}$ stand for the types of its parameters, and $c_{f,1}, \dots, c_{f,m}$ stand for the contexts of its parameters.

Let $C \times [exp] \Pi$ be the total context of all occurrences of the variable x in the expression exp , the expression exp itself being in the context Π .

We have the following set of equations

$$c_{f,j} = \text{TypeToContext}[\text{GenType}[t_{f,j}] (C \times_j [body_f] \{ \{ \} \})]$$

where $C \times [exp] \Pi$ is defined as follows:

$$C \in \text{VName} \rightarrow \text{Exp} \rightarrow \text{Context} \rightarrow \text{Context}$$

$$C \times [x] \Pi = \Pi$$

$$C \times [y] \Pi = \{ \{ \} \}, \quad \text{where } x \neq y.$$

$$C \times [\text{quote } \mathcal{E}] \Pi = \{ \{ \} \}$$

$$C \times [\text{if } exp \text{ then } exp' \text{ else } exp''] \Pi = \\ C \times [exp] \{ \{ \} \} \cup C \times [exp'] \{ \{ \} \} \cup C \times [exp''] \{ \{ \} \}$$

$$C \times [\text{call } f(exp_1, \dots, exp_m)] \Pi = \cup \{ C \times [exp_j] c_{f,j} \}_{j=1, \dots, m}$$

$$C \times [\text{car}(exp)] \Pi = C \times [exp] (\{ \{ \} \} \cup \text{car} * \Pi)$$

$$C \times [\text{cdr}(exp)] \Pi = C \times [exp] (\{ \{ \} \} \cup \text{cdr} * \Pi)$$

$$C \times [\text{cons}(exp', exp'')] \Pi =$$

$$C \times [exp'] (\{ \{ \} \} \cup \Pi / \text{car}) \cup C \times [exp''] (\{ \{ \} \} \cup \Pi / \text{cdr})$$

$$C \times [\text{atom}(exp)] \Pi = C \times [exp] \{ \{ \} \}$$

$$C \times [\text{equal}(exp', exp'')] = C \times [exp'] \{ \{ \} \} \cup C \times [exp''] \{ \{ \} \}$$

where we use the notation $\Pi / \text{car} = \{ \pi \mid [\text{car}]^{\wedge} \pi \in \Pi \}$, $\Pi / \text{cdr} = \{ \pi \mid [\text{cdr}]^{\wedge} \pi \in \Pi \}$.

We assume the set of contexts to be equipped with natural partial ordering, $\Pi' \leq \Pi''$ being equivalent to $\Pi' \subseteq \Pi''$. The functions TypeToContext , GenType , and C are monotonic with respect to contexts, therefore the minimal fixed point for the above system of equations does exist.

Moreover, since $c_{f,j} \subseteq \text{TypeToContext}[t_{f,j}]$, there exist only a finite number of contexts that can be taken as value by $c_{f,j}$, hence the minimal fixed point can be found by a finite number of iterations.

The context analysis above resembles, in some respects, the "neighborhood analysis" as used in the Supercompiler [Turchin 86], [Turchin 88].

14. PRACTICAL IMPLEMENTATION OF THE CONTEXT ANALYSIS

Some programming tricks have prove to be useful for implementing the above backward analysis.

First, what we really use in splitting parameters are types generalized with respect to contexts, rather than contexts themselves. Thus, instead of computing $c_{f,j}$, we can compute the type

$t'_{f,j} = \text{GenType}[t_{f,j}] \ c_{f,j}$.

Second, since $t_{f,j} \leq t'_{f,j}$, we can replace $t_{f,j}$ and $t'_{f,j}$ with a single marked type $mt_{f,j}$ having the syntax

```
mt ∈ MType      marked types
mt ::= any | atom(A) | cons(t',t'') | cons!(mt',mt'') | 1
```

where each marked cons! "belongs" both to $t'_{f,j}$ and $t_{f,j}$, whereas each unmarked cons "belongs" only to $t_{f,j}$, the corresponding place in $t'_{f,j}$ being *any*.

The context $c_{f,j}$ can be extracted from $mt_{f,j}$ directly, without finding $t'_{f,j}$, by means of the function Retrieve! .

```
Retrieve! ∈ MType → Context
```

```
Retrieve![any] = {[ ]}
```

```
Retrieve![atom(A)] = {[ ]}
```

```
Retrieve![cons(t',t'')] = {[ ]}
```

```
Retrieve![cons!(mt',mt'')] = {[ ]} ∪ car*Retrieve![mt'] ∪ cdr*Retrieve![mt'']
```

Next improvement concerns the representation of contexts. Being sets of paths, contexts are difficult to deal with directly, but we can replace contexts with their representations having the syntax

```
crep ∈ ContextRep
```

```
crep ::= car(crep) | cdr(crep) | mtype(mt)
```

Given a context's representation, we can reconstruct the context by the function Retrieve .

```
Retrieve ∈ ContextRep
```

```
Retrieve[car(crep)] = {[ ]} ∪ car*Retrieve[crep]
```

```
Retrieve[cdr(crep)] = {[ ]} ∪ cdr*Retrieve[crep]
```

```
Retrieve[mtype(mt)] = Retrieve![mt]
```

All functions the access path analysis involves can be modified so that they will deal with the representation of contexts, rather than with the contexts themselves.

15. GENERALIZATIONS

The first obvious generalization concerns splitting the results of functions. In the language Mixwell each function produces one and only one result, for which reason a defined function call $\text{call } f(\text{exp}_1, \dots, \text{exp}_m)$ cannot be split and, therefore, has to be assigned the type *any*. The language, however, can be extended, so that a function can produce several results, and this device allows the results of a function to be split

without splitting the function's definition. A version of the arity raiser with this extension has been implemented by Ruten F.Gurin.

Another possible extension is to make an arity raiser deal with data structures that are more complicated than Lisp S-expressions are. For example, in the case of the languages Refal [Turchin 79], [Turchin 86] and RL [Romanenko 88], the data are arbitrary trees, rather than binary trees, which was taken into account in the arity raiser described in [Romanenko 88].

CONCLUSIONS

In order for the results produced by variable splitting to be reasonable, we need information obtained by two preliminary global analyses of the program. The first, forward, analysis tells us whether the splitting is feasible, whereas the second, backward, analysis tells us whether the splitting is useful.

The information obtained is used to avoid introducing new selectors into the program as well as code duplication, and makes it possible to avoid useless variable splitting that does not cause some selectors in the program to be eliminated.

The experiments made by the author have shown that introducing an arity raiser as a separate phase into a specializer enhances the structure of residual programs generated without affecting the other phases of the specializer. The structure of the specializer, thus, can be kept natural and understandable.

REFERENCES

- [Barzdin 88] G.Barzdin. Mixed Computation and Compiler Basis. In D.Bjorner, A.P.Ershov and N.D.Jones, editors, *Partial Evaluation and Mixed Computation*, pages 15-26, North-Holland, 1988.
- [Beckman 76] L.Beckman, A.Haraldson, O.Oskarsson, E.Sandewall. A Partial Evaluator, and Its Use as a Programming Tool. *Artificial Intelligence*, 7(4):319-357, 1976.
- [Bulyonkov 84] M.A.Bulyonkov. Polyvariant Mixed Computation for Analyzer Programs. *Acta Informatica*, 21:473-484, 1984.
- [Burstall 77] R.M.Burstall and J.Darlington. A Transformation System for

Developing Recursive Programs. *Journal of the ACM*, 24(1):44-67, 1977.

- [Dixon 71] J.Dixon. The Specializer, a Method of Automatically Writing Computer Programs. Division of Computer Research and Technology, National Institute of Health, Bethesda, Maryland, 1971.
- [Ershov 78] On the Essence of Compilation. In E.J.Neuhold, editor, *Formal Description of Programming Concepts*, pages 391-420, North-Holland, 1978.
- [Ershov 81] A.P.Ershov. The Transformational Machine: Theme and Variations. In J.Grushka and M.Chytil, editors, *Mathematical Foundations of Computer Science*, Štrbské Pleso, Czechoslovakia, pages 16-32, Lecture Notes in Computer Science, Vol.118, Springer-Verlag, 1981.
- [Futamura 71] Partial Evaluation of Computation Process - An Approach to a Compiler-Compiler. *Systems, Computers, Controls*, 2(5):45-50, 1971.
- [Hughes 88] J.Hughes. Backward Analysis of Functional Programs. In D.Bjorner, A.P.Ershov and N.D.Jones, editors, *Partial Evaluation and Mixed Computation*, pages 187-208, North-Holland, 1988.
- [Jones 85] N.D.Jones, P.Sestoft and H.Sondergaard. An Experiment in Partial Evaluation: The Generation of a Compiler Generator. In J.-P.Jouannaud, editor, *Rewriting Techniques and Applications, Dijon, France*, pages 124-140, Lecture Notes in Computer Science, Vol.202, Springer-Verlag, 1985.
- [Jones 86] N.D.Jones and A.Mycroft. Data Flow Analysis of Applicative Programs Using Minimal Function Graphs. In *Thirteens ACM Symposium on Principles of Programming Languages, St.Petersburg, Florida*, pages 296-306, ACM, 1986.
- [Jones 88] Automatic Program Specialization: A Re-Examination from Basic Principles. In D.Bjorner, A.P.Ershov and N.D.Jones, editors, *Partial Evaluation and Mixed Computation*, pages 225-282, North-Holland, 1988.
- [Mogensen 88] T.Mogensen. Partially Static Structures in a Self-Applicable Partial Evaluator. In D.Bjorner, A.P.Ershov and N.D.Jones, editors, *Partial Evaluation and Mixed Computation*, pages 325-347, North-Holland, 1988.
- [Ostrowski 88] B.N.Ostrowski. Implementation of Controlled Mixed Computation in System for Automatic Development of Language-Oriented Parsers. In D.Bjorner, A.P.Ershov and N.D.Jones, editors, *Partial*

Evaluation and Mixed Computation, pages 385-403, North-Holland, 1988.

- [Romanenko 88] S.A.Romanenko. A Compiler Generator Produced by a Self-Applicable Specializer Can Have a Surprisingly Natural and Understandable Structure. In D.Bjorner, A.P.Ershov and N.D.Jones, editors, *Partial Evaluation and Mixed Computation*, pages 445-463, North-Holland, 1988.
- [Sestoft 86] The Structure of a Self-Applicable Partial Evaluator. In H.Ganzinger and N.D.Jones, editors, *Programs as Data Objects, Copenhagen, Denmark, 1985*, pages 236-256, Lecture Notes in Computer Science, Vol. 217, Springer-Verlag, 1986.
- [Schmidt 86] D.A.Schmidt. *Denotational Semantics*. Allyn and Bacon, Boston, 1986.
- [Sestoft 88] P.Sestoft. Automatic Call Unfolding in a Partial Evaluator. In D.Bjorner, A.P.Ershov and N.D.Jones, editors, *Partial Evaluation and Mixed Computation*, pages 485-506, North-Holland, 1988.
- [Turchin 72] V.F.Turchin. Equivalent Transformation of Recursive Functions Defined in Refal. In *Teoriya Yazykov i Metody Programirovaniya. Trudy Simposiuma*, pages 31-42, Alushta-Kiev, 1972 (in Russian).
- [Turchin 79] V.F.Turchin. A Supercompiler System Based on the Language Refal. *SIGPLAN Notices*, 14(2):46-54, February 1979.
- [Turchin 82] V.F.Turchin, R.M.Nirenberg and D.V.Turchin. Experiments with a Supercompiler. In *1982 ACM Symposium on Lisp and Functional Programming, Pittsburgh, Pennsylvania*, pages 47-55, ACM, 1982.
- [Turchin 86] V.F.Turchin. The Concept of a Supercompiler. *ACM Transactions on Programming Languages and Systems*, 8(3):292-325, July 1986.
- [Turchin 88] V.F.Turchin. The Algorithm of Generalization in the Supercompiler. In D.Bjorner, A.P.Ershov and N.D.Jones, editors, *Partial Evaluation and Mixed Computation*, pages 531-549, North-Holland, 1988.
- [Wadler 88] P.Wadler. Deforestation: Transforming Programs to Eliminate Trees. In *European Symposium on Programming, Lecture Notes in Computer Science*, Springer-Verlag, 1988.