

Efficient Compilation of Lazy Evaluation

Thomas Johnsson

Programming Methodology Group
Department of Computer Science
Chalmers University of Technology
S-412 96 Göteborg, Sweden

Abstract

This paper describes the principles underlying an efficient implementation of a lazy functional language, compiling to code for ordinary computers. It is based on combinator-like graph reduction: the user defined functions are used as rewrite rules in the graph. Each function is compiled into an instruction sequence for an abstract graph reduction machine, called the G-machine, the code reduces a function application graph to its value. The G-machine instructions are then translated into target code. Speed improvements by almost two orders of magnitude over previous lazy evaluators have been measured; we provide some performance figures.

1. Background

Functional programming is emerging as an alternative to the conventional imperative style of programming [Land66, Back78]. Lazy evaluation (call by need, normal order evaluation) has been proposed as a method for executing functional programs, the advantages being, among others, that unbound data structures, e.g. infinite lists, can easily be handled, and further that it makes interactive input/output possible in functional programs [Frie76]. Though functional programming languages have many pleasing properties, an obstacle to their wider use has been the lack of efficient implementations.

Our work is based on Turner's combinator approach [Turn79], where programs are transformed into expressions containing the combinators S, K, I etc from combinatory logic, thus removing all variables from the program. A combinator expression is evaluated in the 'SKI-machine' using normal order graph reduction. A problem with combinators is that each combinator defines rather a small interpretative step, and combinator expressions have a tendency to become very cumbersome for non-trivial programs.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

©1984 ACM 0-89791-139-3/84/0600/0058\$00.75

Our lazy evaluation method is similar to the combinator reduction regime, but instead of using a standard, fixed set of combinators, each user defined function is used as a 'combinator', i.e., a rewrite rule for the graph. Functions are compiled into code sequences for an abstract graph reduction machine, called the G-machine, with instructions that explicitly construct and manipulate expression graphs to reduce expressions to their values; both shared and cyclic graphs can be directly constructed. Target code generation for ordinary computers from the G-machine code is rather straightforward. One might say that the compiler constructs a specialised, machine-language coded combinator interpreter from each program.

2. Graph reduction

In our graph reduction approach a program is an expression whose value will appear as, in general, a stream of basic values (integers, booleans etc) on the output file. Expressions are evaluated using normal order graph reduction, and is carried out by performing transformations on the graph to reduce it to its **canonical form**. A canonical form is an expression which cannot be further reduced on the outermost level (even though subexpressions may be further reducible). In this paper canonical forms are integer and boolean constants, list expressions $e_1.e_2$ with arbitrary expressions e_1 and e_2 , and function applications $f e_1 \dots e_m$ where f is a function that takes more than m curried arguments; a reduction of an application can take place only if all curried arguments to the functions are present. Thus in general for an expression to become completely reduced, subexpressions must also be reduced to canonical form, for instance the elements of a list. Evaluation of a function application amounts to using the corresponding function definition as a graph rewrite rule, repeatedly rewriting the application graph to an instance of the right hand side of the function definition, with arguments substituted for formal parameters, until having reached a canonical form.

For illustration, consider the following functional program, its value being the infinite list of natural numbers.

```
letrec from n = n.from(succ n) in from 0
```

Succ is the predefined successor function, '.' is the infix list construction operator and juxtaposition denotes function application. Graph reduction of this program is shown in figure 1. In the figures function application is denoted by @.

The start expression 1(a) is transformed to 1(b) using the rewrite rule for the function 'from' as defined above, with a pointer to the integer 0 substituted for the parameter 'n'. In 1(b) the expression is on canonical form, and so is also its head part 0. The head value can now be output and dropped from the graph, 1(c). Again the rewrite rule for 'from' is applied to the graph, 1(d), and is now on the form e.e', which is canonical. The next step is to reduce the head part 'succ 0' to its canonical form using the rewrite rule for 'succ', 1(e). The resulting integer 1 is then written on output and dropped from the graph.

The execution continues in this way ad infinitum. Note that the shared expression 'succ 0' has been replaced in 1(e) by its value. In general expression graphs are evaluated, i.e. reduced to their canonical forms, at most once and all expressions that share a particular subexpression benefit from the evaluation (call by need).

3. An introductory example of G-machine execution

In our graph reduction scheme each function definition is compiled into a sequence of G-machine instructions. Each graph rewrite, according to a function definition, is carried out by executing the code for that function. We here illustrate execution of the G-machine with the reduction step (c)-(d) from figure 1. The G-machine state transitions are shown in figure 2.

Before the start of the reduction a pointer to the expression graph is at the top of a pointer stack, figure 2(a) (the stack grows downwards). Reduction is started by execution of the G-machine instruction EVAL, in this case by the print mechanism.

EVAL causes a new stack frame to be created with the previously topmost pointer as its single entry, saving the old stack on another stack called the dump (not shown in figure 2), then an unwind state pushes pointers to the application nodes of the left 'backbone' until having reached a function node, (a) - (c). The stack is then rearranged so that the topmost pointer of the stack points to the argument of 'from', the second pointer from the top is left untouched and will thus point to the apply node which is to be updated by the code with the result of the application. The G-machine now starts to execute the code for 'from', which is (see section 5 and table 3 how we obtain this)

```
from: PUSH 0; PUSHFUN from;
      PUSHFUN succ; PUSH 3; MKAP;
      MKAP; CONS; UPDATE 2; RET 1.
```

Except for the last two instructions, this instruction sequence is essentially a postfix representation of the right hand side of 'from'. The PUSH m instruction pushes the m'th pointer of the stack relative to the top and starting with 0; note that different offsets have to be used to push pointers to the formal parameter n, depending on the current depth of the stack (the reason for this is explained in section 5.5). The PUSHFUN succ instruction pushes a pointer to a 'succ' function node. MKAP constructs an application node with the to topmost as subparts; similarly for CONS. After having constructed the graph for the right hand side of the definition of 'from', figure 2(k), the cons node is copied onto the result apply node by the UPDATE 2 instruction, having thereby transformed 'from(succ 0)' to '(succ 0).from(succ 0))' in the graph, which is a canonical expression. The RET 1 instruction pops one element from the stack, and since the top graph is now on canonical form, the old stack is restored from the dump and control is returned to the instruction following EVAL. In general, had the top graph been not on canonical form but an application node or a function node, instead of

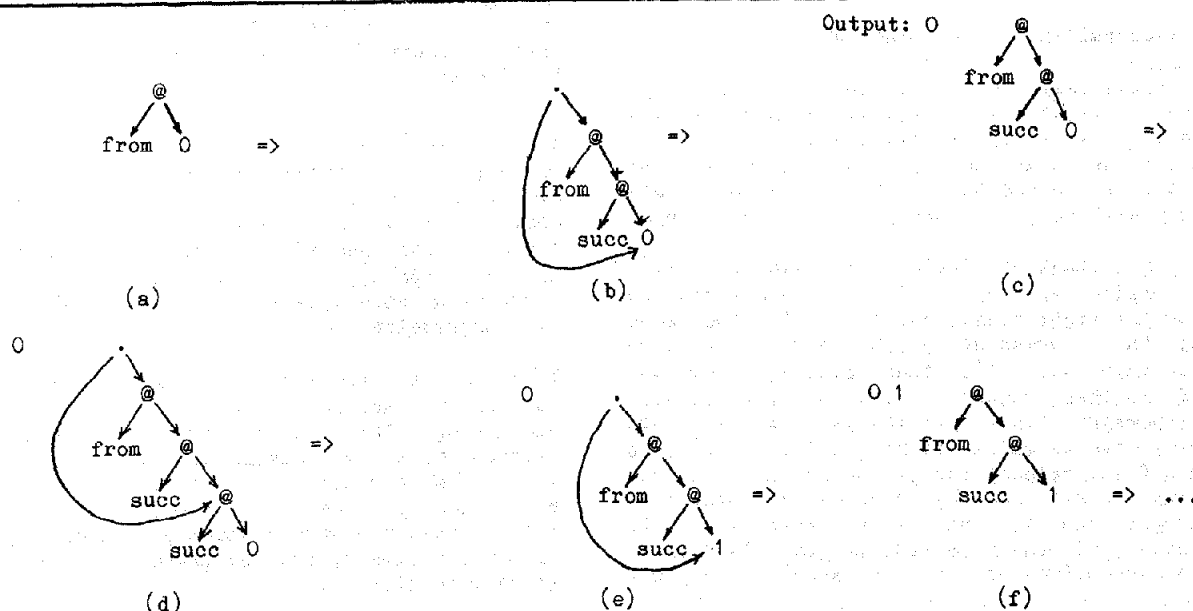


Figure 1. Graph reduction of 'from 0'. The output is shown to the left of each graph.

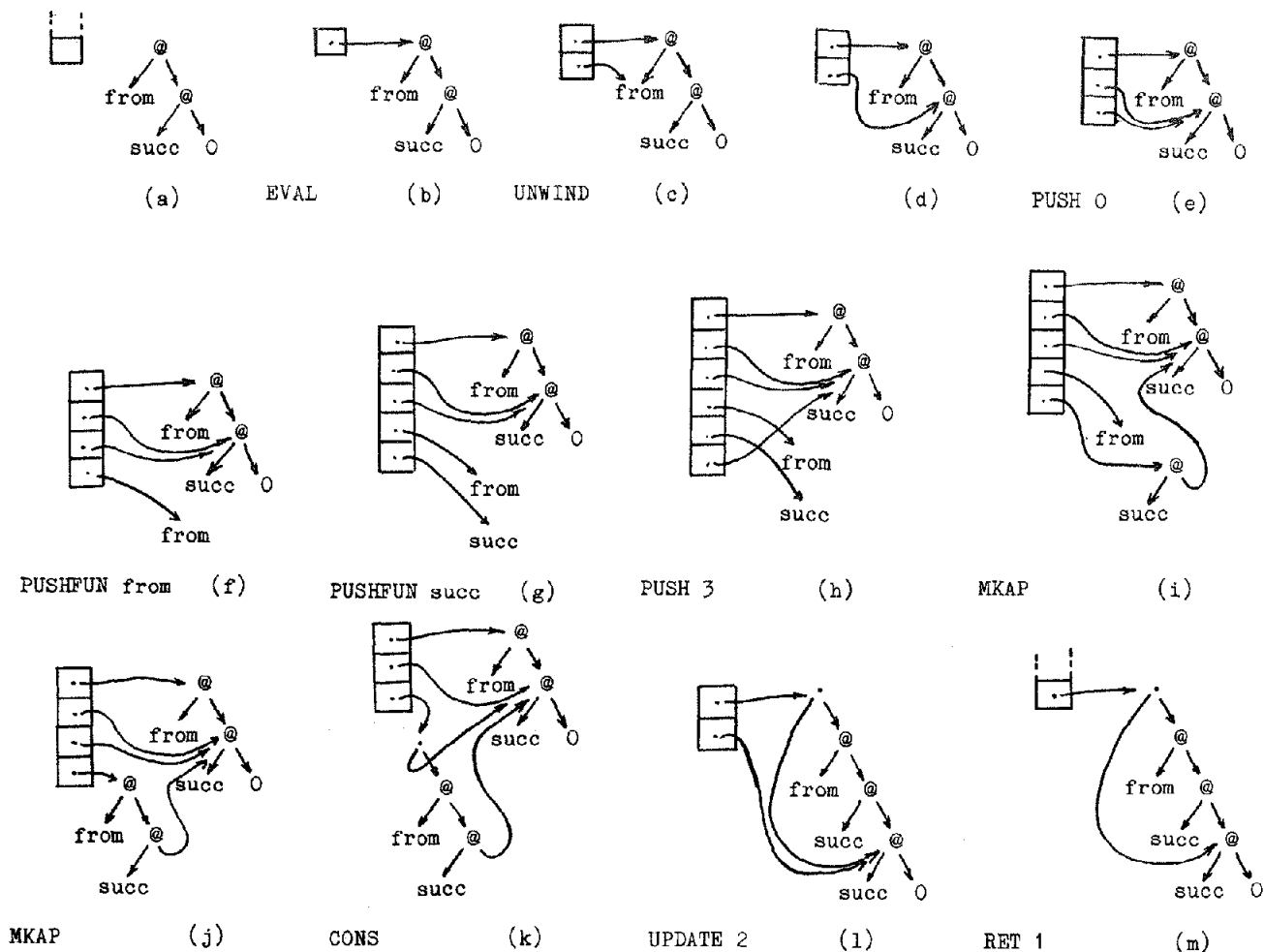


Figure 2. G-machine reduction of 'from (succ 0)'.

restoring the old stack and returning the G-machine would have reentered the unwind state to continue the reduction of the new expression graph.

4. Short-circuiting graph reduction

We have previously indicated that we do graph reduction by repeatedly rewriting the graph to the right hand sides of functions. Indeed we can use G-machine code that does precisely this; in most cases, however, we can take considerable shortcuts and do away with many intermediate graph rewritings.

Consider the function definition 'succ $n = n+1$ '. If we compile it into code that constructs the graph for the right hand side, 'add $n\ 1$ ', then when executed the expression graph 'succ e ' will be rewritten into 'add $e\ 1$ ', thus leaving over the task of further reduction to 'add', which will reduce the expression to its integer value. Much efficiency can be gained if we compile 'succ' into code that first reduces its parameter n , computes the value of $n+1$, and then remakes the apply node to a integer node with this value. This avoids the construction of the intermediate graph 'add $e\ 1$ '. A code sequence for the function 'succ' is accordingly

```
succ: PUSH 0; EVAL; GET; PUSHBASIC 1; ADD;
      MKINT; UPDATE 2; RET 1.
```

The execution of this code sequence is shown in figure 3. The addition is done on a separate stack for basic values, called V, with instructions MKINT and GET for transferring values to and from the graph. PUSHBASIC pushes a basic value constant on the V stack.

Similar reasoning can be applied to all other predefined primitive functions; if the right hand side is an if-expression, for example, then the code would do the following: compute the value of the condition, and if true the proper apply node is to be updated with the value of the then-expression, else updated with the value of the else-expression.

This line of reasoning is systematized in the next section by having different compilation schemes, one giving code that computes the value of an expression, and one giving code that constructs the graph of an expression. This more direct method is significantly faster; in our compiler implementation we have measured a speedup of about a factor of ten for some typical programs, compared to naive graph reduction.

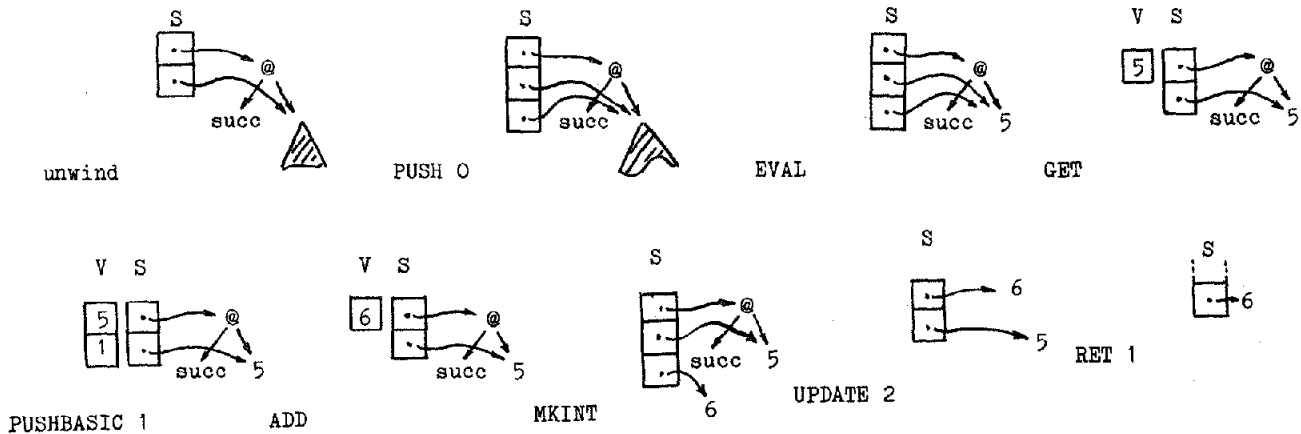


Figure 3. Shortcut evaluation of function succ.

5. Technical details of the abstract machine and compiler

In this section we give a complete set of compilation rules for a simple functional language, compiling to G-machine code. We also give an abstract description of the G-machine, describing the effects of the G-machine instructions on a machine state.

5.1 Source Language

A program in the language described here consists of a set of recursive functions and an expression whose value is the value of the program, as summarised in table 1. Normal order evaluation is assumed. Each function f_i takes n_i curried arguments and the free variables of e_i are in the set $\{x_1 \dots x_{n_i}\}$. Operators $+$, $-$ etc are viewed as syntactic sugar for applications to predefined functions add, sub etc, of which we deal with the ones given in table 2.

TABLE 1. Syntax of programs

```

program ::=
  f1 x1 ... xn1 = e1    (function definitions)
  ...
  fm x1 ... xnm = em
  e0                        (the value of the program)

e ::= identifiers | constants | e e |
    let x1=e1 and ... and xm=em in e |
    (multiple simultaneous local definitions)
    letrec x1=e1 and ... and xm=em in e
    (multiple simultaneous local recursive definitions)
  
```

TABLE 2. Predefined functions

add sub mul div	(binary arithmetic operators)
neg	(unary negation)
lt le eq ne ge gt	(binary relational operators)
and or	(conditional and, or)
not	(logical negation)
cons	(binary list construction)
hd, tl	(unary head and tail of a list)
null	(unary test on empty list)
if	(ternary if-then-else)

Note that there is no lambda expression in the syntax of expressions, functions are defined only globally. Functional programs with local function definitions and lambda expressions with free variables can be transformed into the form above, using 'super combinators' [Hugh82]; an algorithm to the same effect is used in our compiler implementation, however, the program resulting from our transformation does not exhibit 'fully laziness', as is the main issue in Hughes' work.

5.2 Compilation rules

The abstract compiler given in table 3 is subdivided into 4 compilation schemes:

$F[f x_1 \dots x_{n_i} = e]$ gives the code for a function which reduces the graph of an application to canonical form.

$C[e] r n$ gives code that constructs the graph of e and leaves a pointer to the result on the top of the stack.

$E[e] r n$ gives code that computes the value, i.e. canonical form, of e and leaves a pointer to the value on the top of the stack. It yields the same result as $C[e] r n$ followed by an EVAL instruction and embodies the short-circuiting described in section 4.

$B[e] r n$ computes the basic value of e and leaves the result on the basic value stack V , yielding the same result as $E[e] r n$ followed by a GET instruction. The idea behind B is to avoid construction of a new node for each intermediate result in an arithmetic or logical expression. The value is transferred to the graph only when the entire expression has been evaluated, by the MKINT or MKBOOL instruction.

In addition, there are 3 help-functions used for local definitions: Xr returns a pair of the extended environment and the new stack depth, $Clet$ and $Cletrec$ gives code to extend the stack with local definitions. In the translation schemes r is a mapping from identifiers of parameters to their

TABLE 3. Compilation rules

$F[f, x_1, \dots, x_m = e] = E[e] \text{ r } (m+1); \text{UPDATE } (m+1); \text{RET } m$, where $r = [x_1=m+1, x_2=m, \dots, x_m=2]$

Scheme E: Evaluate

1. $E[i] \text{ r } n = \text{PUSHINT } i$
2. $E[b] \text{ r } n = \text{PUSHBOOL } b$
3. $E[\text{nil}] \text{ r } n = \text{PUSHNIL}$
4. $E[x] \text{ r } n = \text{PUSH } (n - r \text{ x}); \text{EVAL}$
5. $E[f] \text{ r } n = \text{PUSHFUN } f$
6. $E[\text{add } e_1 e_2] \text{ r } n = E[\text{add } e_1 e_2] \text{ r } n; \text{MKINT, and similarly for sub, mul, div}$
7. $E[\text{neg } e] \text{ r } n = E[\text{neg } e] \text{ r } n; \text{MKINT}$
8. $E[\text{eq } e_1 e_2] \text{ r } n = E[\text{eq } e_1 e_2] \text{ r } n; \text{MKBOOL, and similarly for lt, gt, ne, ge, le}$
9. $E[\text{not } e] \text{ r } n = E[\text{not } e] \text{ r } n; \text{MKBOOL}$
10. $E[\text{and } e_1 e_2] \text{ r } n = E[\text{if } e_1 \text{ false}] \text{ r } n$
11. $E[\text{or } e_1 e_2] \text{ r } n = E[\text{if } e_1 \text{ true } e_2] \text{ r } n$
12. $E[\text{cons } e_1 e_2] \text{ r } n = C[e_1] \text{ r } n; C[e_2] \text{ r } (n+1); \text{CONS}$
13. $E[\text{null } e] \text{ r } n = E[e] \text{ r } n; \text{NULL}; \text{MKBOOL}$
14. $E[\text{hd } e] \text{ r } n = E[e] \text{ r } n; \text{HD}; \text{EVAL, similarly for tl}$
15. $E[\text{if } e_1 e_2 e_3] \text{ r } n = E[e_1] \text{ r } n; \text{JFALSE } l_1; E[e_2] \text{ r } n; \text{JMP } l_2; \text{LABEL } l_1; E[e_3] \text{ r } n; \text{LABEL } l_2$
where l_1 and l_2 are new unique labels
16. $E[\text{let } d \text{ in } e] \text{ r } n = \text{Clet}[d] \text{ r } n; E[e] \text{ r } n'; \text{SLIDE } (n'-n) \text{ where } (r', n') = \text{Xr}[d] \text{ r } n$
17. $E[\text{letrec } d \text{ in } e] \text{ r } n = \text{Cletrec}[d] \text{ r } n'; E[e] \text{ r } n'; \text{SLIDE } (n'-n) \text{ where } (r', n') = \text{Xr}[d] \text{ r } n$
18. $E[e] \text{ r } n = C[e] \text{ r } n; \text{EVAL otherwise}$

Scheme B: Evaluate basic value

1. $B[i] \text{ r } n = \text{PUSHBASIC } i$
2. $B[b] \text{ r } n = \text{PUSHBASIC } b$
3. $B[\text{add } e_1 e_2] \text{ r } n = B[e_1] \text{ r } n; B[e_2] \text{ r } (n+1); \text{ADD; similarly for sub, mul, div, eq, ne, lt, gt, ge, le.}$
4. $B[\text{neg } e] \text{ r } n = B[e] \text{ r } n; \text{NEG}$
5. $B[\text{not } e] \text{ r } n = B[e] \text{ r } n; \text{NOT}$
6. $B[\text{null } e] \text{ r } n = E[e] \text{ r } n; \text{NULL}$
7. $B[\text{if } e_1 e_2 e_3] \text{ r } n = B[e_1] \text{ r } n; \text{JFALSE } l_1; B[e_2] \text{ r } n; \text{JMP } l_2; \text{LABEL } l_1; B[e_3] \text{ r } n; \text{LABEL } l_2$
where l_1 and l_2 are new unique labels
8. $B[\text{let } d \text{ in } e] \text{ r } n = \text{Clet}[d] \text{ r } n; B[e] \text{ r } n'; \text{POP } (n'-n) \text{ where } (r', n') = \text{Xr}[d] \text{ r } n$
9. $B[\text{letrec } d \text{ in } e] \text{ r } n = \text{Cletrec}[d] \text{ r } n'; B[e] \text{ r } n'; \text{POP } (n'-n) \text{ where } (r', n') = \text{Xr}[d] \text{ r } n$
10. $B[e] \text{ r } n = E[e] \text{ r } n; \text{GET otherwise}$

Scheme C: Construct graph

1. $C[i] \text{ r } n = \text{PUSHINT } i$
2. $C[b] \text{ r } n = \text{PUSHBOOL } b$
3. $C[\text{nil}] \text{ r } n = \text{PUSHNIL}$
4. $C[f] \text{ r } n = \text{PUSHFUN } f$
5. $C[x] \text{ r } n = \text{PUSH } (n - r \text{ x})$
6. $C[\text{cons } e_1 e_2] \text{ r } n = C[e_1] \text{ r } n; C[e_2] \text{ r } (n+1); \text{CONS}$
7. $C[e_1 e_2] \text{ r } n = C[e_1] \text{ r } n; C[e_2] \text{ r } (n+1); \text{MKAP, if not matched above}$
8. $C[\text{let } d \text{ in } e] \text{ r } n = \text{Clet}[d] \text{ r } n; C[e] \text{ r } n'; \text{SLIDE } (n'-n) \text{ where } (r', n') = \text{Xr}[d] \text{ r } n$
9. $C[\text{letrec } d \text{ in } e] \text{ r } n = \text{Cletrec}[d] \text{ r } n'; C[e] \text{ r } n'; \text{SLIDE } (n'-n) \text{ where } (r', n') = \text{Xr}[d] \text{ r } n$

Miscellaneous schemes for local definitions

$\text{Xr}[v_1=e_1 \text{ and } \dots v_i=e_i \text{ and } \dots v_m=e_m] \text{ r } n = (r[v_1=n+1, \dots v_i=n+i, \dots v_m=n+m], n+m)$
 $\text{Clet}[v_1=e_1 \text{ and } \dots v_i=e_i \text{ and } \dots v_m=e_m] \text{ r } n = C[e_1] \text{ r } n; \dots C[e_i] \text{ r } (n+i-1); \dots C[e_m] \text{ r } (n+m-1)$
 $\text{Cletrec}[v_1=e_1 \text{ and } \dots v_i=e_i \text{ and } \dots v_m=e_m] \text{ r } n = \text{ALLOC } m; C[e_1] \text{ r } (n+m); \text{UPDATE } m; \dots$
 $C[e_i] \text{ r } (n+m); \text{UPDATE } (m+1-i); \dots C[e_m] \text{ r } (n+m); \text{UPDATE } 1$

location on the stack, and n is the current depth of the stack. Below we show compilation of the function $f\ x = x.(f\ x)$.

```
F[f x = cons x (f x)] =
E[cons x (f x)] [x=2] 2; UPDATE 2; RET 1 =
C[x] [x=2] 2; C[f x] [x=2] 3;
CONS; UPDATE 2; RET 1 =
PUSH 0; C[f] [x=2] 3; C[x] [x=2] 4; MKAP;
CONS; UPDATE 2; RET 1 =
PUSH 0; PUSHFUN f; PUSH 2; MKAP;
CONS; UPDATE 2; RET 1.
```

5.3 The abstract machine

A state in the abstract G-machine is a 7-tuple $\langle O, C, S, V, G, E, D \rangle$ where

- O is the output produced so far, as shown in the example in figure 1. It consists of a sequence of integers and booleans. In an actual implementation O is printed on standard output.
- C is the G-code sequence currently being executed.
- S is a stack of node names, i.e. pointers into the graph.
- V is a stack of basic values, i.e. integers and booleans on which the arithmetic and logical operations are performed, as shown in section 4.
- G is the graph: a mapping from node names to nodes. We have nodes of the following types:
 - INT i integer nodes,
 - BOOL b boolean nodes,
 - NIL empty list nodes,
 - CONS $n_1\ n_2$ list nodes, where n_1 is a pointer to the head graph and n_2 is a pointer to the tail graph,
 - AP $n_1\ n_2$ application nodes, where n_1 is a pointer to the function graph and n_2 is a pointer to the argument graph,
 - FUN f a node with a reference to the compiled function f ,
 - HOLE a node which is to be filled in with another value later; it is used while constructing cyclic graphs for letrec expressions.
- E is a global environment, which is a mapping from function names to pairs consisting of the number of curried arguments of the function, and its code sequence. E corresponds to the code segment in conventional machines and is constant throughout the execution of the program.
- D is a dump used for recursive calls to EVAL: a stack of pairs consisting of
 - a stack of node names: S before EVAL,
 - a G-code sequence: C before EVAL.

Table 4 summarises the state transition rules for the G-machine instructions used in the compilation rules given in table 3.

In a G-machine state, $()$ denotes an empty stack or an empty code sequence. The semicolon appends values onto an output sequence. Period is used as infix cons for instruction sequences and push for stacks. Updating of the graph is written as e.g. $G[n = \text{INT } i]$. If there is a node named n previ-

ously in G , then then node n is updated with a new value, otherwise a new node is created. This notation is also used in pattern matching situations, for instance state transition rule 1 is applicable if the top of the stack points to an integer node. For instructions with parameters, e.g. PUSH $m.c$ binds as (PUSH m).c. A node name that occurs only in the right hand side of a transition rule is considered to be new and unique, e.g. n' in transition rule 12. G-machine states that do not match any rule are considered to be run time errors.

The definition of the G-machine has certain similarities with the definition of the SECD machine [Land64], new in our model is that we describe how we do lazy output, and handle updating and sharing in a graph, in the framework of the abstract machine.

5.4 Initial and final state of the machine

The initial configuration of the machine for a given program is shown in figure 4. The machine starts with an empty output, a code sequence c_0 for evaluating and printing the start expression e_0 , an empty pointer stack and an empty basic value stack, an empty graph, an environment E_0 containing the compiled code for the functions together with their arity, and an empty dump.

$\langle (), c_0, (), (), \{\}, E_0, () \rangle$
 where $c_0 = E[e_0]_{r_0 0}$; PRINT
 and $E_0 = \{ f_0 : (n_1, F[f_0 x_1 \dots x_{n_1} = e_0]), \dots$
 $f_m : (n_m, F[f_m x_1 \dots x_{n_m} = e_m])$
 $\text{add} : (2, F[p\ x\ y = \text{add } x\ y])$
 $\text{sub} : (2, F[p'\ x\ y = \text{sub } x\ y])$
 \dots
 $\}$

Figure 4. Initial state of the G-machine.

Since the operators $+$, $-$ etc are represented with applications to predefined functions add, sub etc in unevaluated expression graphs, the code for these functions must also be present in E_0 . The machine stops when the state $\langle 0, (), (), (), G, E, () \rangle$ has been reached.

5.5 The evaluation mechanism

The evaluation of the program is driven by PRINT, which in case of a list starts the evaluation of the head and the tail part of the list, see transition rules 1-4. Only the leaves of the printed data structure appears on the output, for instance the list $(2.3.\text{nil}).(5.\text{nil}).\text{nil}$ gives the output sequence 2 3 5.

The EVAL instruction reduces the graph pointed to by the pointer at the top of the stack to canonical form. If the top of stack is an apply node, transition rule 5, the rest of the code sequence and the stack except for the top element is pushed onto the dump, and the unwind state is entered, following the function parts of apply nodes, pushing the function pointers on the way (transition rule 7). When a function node has been reached, and the stack is deep enough to contain all curried arguments to the function, rule 8, the stack is arranged according to figure 5. The top m elements

TABLE 4. State transition rules for G-machine instructions

1. $\langle o, \text{PRINT.c, n.s, v, G[n=INT i], E, D} \rangle \Rightarrow \langle o; i, c, s, v, G[n=INT i], E, D \rangle$
2. $\langle o, \text{PRINT.c, n.s, v, G[n=BOOL b], E, D} \rangle \Rightarrow \langle o; b, c, s, v, G[n=BOOL b], E, D \rangle$
3. $\langle o, \text{PRINT.c, n.s, v, G[n=CONS } n_1 n_2], E, D \rangle \Rightarrow$
 $\langle o, \text{EVAL.PRINT.EVAL.PRINT.c, } n_1.n_2.s, v, G[n=CONS } n_1 n_2], E, D \rangle$
4. $\langle o, \text{PRINT.c, n.s, v, G[n=NIL], E, D} \rangle \Rightarrow \langle o, c, s, v, G[n=NIL], E, D \rangle$
5. $\langle o, \text{EVAL.c, n.s, v, G[n=AP } n_1 n_2], E, D \rangle \Rightarrow \langle o, \text{UNWIND.(), n.(), v, G[n=AP } n_1 n_2], E, (c,s).D \rangle$
6. $\langle o, \text{EVAL.c, n.s, v, G[n=INT i], E, D} \rangle \Rightarrow \langle o, c, n.s, v, G[n=INT i], E, D \rangle,$
 similarly for nodes BOOL b, NIL, CONS $n_1 n_2$ and FUN f.
7. $\langle o, \text{UNWIND.(), n.s, v, G[n=AP } n_1 n_2], E, D \rangle \Rightarrow \langle o, \text{UNWIND.(), } n_1.n.s, G[n=AP } n_1 n_2], E, D \rangle$
8. $\langle o, \text{UNWIND.(), } n_0.n_1 \dots n_k.s, v, G[n_0=\text{FUN f, } n_1=\text{AP } n_1' n_1'', \dots, n_k=\text{AP } n_k' n_k''], E[f=(k,c)], D \rangle \Rightarrow$
 $\langle o, c, n_1' \dots n_k' .s, v, G[n_0=\text{FUN f, } n_1=\text{AP } n_1' n_1'', \dots, n_k=\text{AP } n_k' n_k''], E[f=(k,c')], D \rangle$
9. $\langle o, \text{UNWIND.(), } n_0.n_1 \dots n_k.(), v, G[n_0=\text{FUN f}], E[f=(a,c')], (c',s').D \rangle$ and $k < a \Rightarrow$
 $\langle o, c', n_k.s', v, G[n_0=\text{FUN f}], E[f=(k,c')], D \rangle$
10. $\langle o, \text{RET m.c, v, } n_1 \dots n_m.n.(), G[n=INT i], E, (c',s').D \rangle \Rightarrow$
 $\langle o, c', n.s', v, G[n=INT i], E, D \rangle,$ similarly for nodes BOOL b, NIL and CONS $n_1 n_2$.
11. $\langle o, \text{RET m.c, } n_1 \dots n_m.n.s, v, G[n=AP } n_1 n_2], E, D \rangle \Rightarrow$
 $\langle o, \text{UNWIND.(), n.s, v, G[n=AP } n_1 n_2], E, D \rangle,$ similarly for $n = \text{FUN f}$.
12. $\langle o, \text{PUSHINT i.c, s, v, G, E, D} \rangle \Rightarrow \langle o, c, n'.s, v, G[n'=INT i], E, D \rangle$
13. $\langle o, \text{PUSHBOOL b.c, s, v, G, E, D} \rangle \Rightarrow \langle o, c, n'.s, v, G[n'=BOOL b], E, D \rangle$
14. $\langle o, \text{PUSHNIL.c, s, v, G, E, D} \rangle \Rightarrow \langle o, c, n'.s, v, G[n'=NIL], E, D \rangle$
15. $\langle o, \text{PUSHFUN f.c, s, v, G, E, D} \rangle \Rightarrow \langle o, c, n'.s, v, G[n'=FUN f], E, D \rangle$
16. $\langle o, \text{PUSH m.c, } n_0 \dots n_m.s, v, G, E, D \rangle \Rightarrow \langle o, c, n_0.n_0 \dots n_m.s, v, G, E, D \rangle$
17. $\langle o, \text{MKINT.c, s, i.v, G, E, D} \rangle \Rightarrow \langle o, c, n'.s, v, G[n'=INT i], E, D \rangle$
18. $\langle o, \text{MKBOOL.c, s, b.v, G, E, D} \rangle \Rightarrow \langle o, c, n'.s, v, G[n'=BOOL b], E, D \rangle$
19. $\langle o, \text{MKAP.c, } n_1.n_2.s, G, E, D \rangle \Rightarrow \langle o, c, n'.s, v, G[n'=AP } n_2 n_1], E, D \rangle$
20. $\langle o, \text{CONS.c, } n_1.n_2.s, G, E, D \rangle \Rightarrow \langle o, c, n'.s, v, G[n'=CONS } n_2 n_1], E, D \rangle$
21. $\langle o, \text{ALLOC m.c, s, v, G, E, D} \rangle \Rightarrow \langle o, c, n_1' \dots n_m'.s, v, G[n_1'=\text{HOLE, } \dots, n_m'=\text{HOLE}], E, D \rangle$
22. $\langle o, \text{UPDATE m.c, } n_0 \dots n_m.s, v, G[n_0: N_0, n_m: N_m], E, D \rangle \Rightarrow$
 $\langle o, c, s, v, n_1 \dots n_m.s, G[n_0: N_0, n_m: N_m], E, D \rangle$
23. $\langle o, \text{SLIDE m.c, } n_0 \dots n_m.s, v, G, E, D \rangle \Rightarrow \langle o, c, n_0.s, v, G, E, D \rangle$
24. $\langle o, \text{GET.c, n.s, v, G[n=INT i], E, D} \rangle \Rightarrow \langle o, c, s, i.v, G[n=INT i], E, D \rangle$
25. $\langle o, \text{GET.c, n.s, v, G[n=BOOL b], E, D} \rangle \Rightarrow \langle o, c, s, b.v, G[n=BOOL b], E, D \rangle$
26. $\langle o, \text{PUSHBASIC i.c, s, v, G, E, D} \rangle \Rightarrow \langle o, c, s, i.v, G, E, D \rangle$
27. $\langle o, \text{ADD.c, s, } i_2.i_1.v, G, E, D \rangle \Rightarrow \langle o, c, s, (i_1+i_2).v, G, E, D \rangle,$ similarly for SUB, MUL, DIV, EQ, NE, LT, GT, LE and GE, the last six putting boolean values on V
28. $\langle o, \text{NEG.c, s, i.v, G, E, D} \rangle \Rightarrow \langle o, c, s, (-i).v, G, E, D \rangle$
29. $\langle o, \text{NOT.c, s, b.v, G, E, D} \rangle \Rightarrow \langle o, c, s, (\text{not } b).v, G, E, D \rangle$
30. $\langle o, \text{JFALSE l.c, s, true.v, G, E, D} \rangle \Rightarrow \langle o, c, s, G, E, D \rangle$
31. $\langle o, \text{JFALSE l.c, s, false.v, G, E, D} \rangle \Rightarrow \langle o, \text{JMP l.c, s, v, G, E, D} \rangle$
32. $\langle o, \text{JMP l. } \dots \text{ LABEL l.c, s, v, G, E, D} \rangle \Rightarrow \langle o, c, s, v, G, E, D \rangle$
33. $\langle o, \text{LABEL l.c, s, G, E, D} \rangle \Rightarrow \langle o, c, s, G, E, D \rangle$
34. $\langle o, \text{HD.c, n.s, v, G[n=CONS } n_1 n_2], E, D \rangle \Rightarrow \langle o, c, n_1.s, v, G[n=CONS } n_1 n_2], E, D \rangle,$
 similarly for TL
35. $\langle o, \text{NULL.c, n.s, v, G[n=CONS } n_1 n_2], E, D \rangle \Rightarrow \langle o, c, s, \text{false.v, G[n=CONS } n_1 n_2], E, D \rangle$
36. $\langle o, \text{NULL.c, n.s, v, G[n=NIL], E, D} \rangle \Rightarrow \langle o, c, s, \text{true.v, G[n=NIL], E, D} \rangle$

of the stack now points to the m carried arguments of the function, and below them there is a pointer to the apply node which is to be updated with the value of the application. The reason for remaking the stack in this manner is firstly to make the function arguments easily accessible, and secondly to access function arguments and local variables introduced by let and letrec expressions uniformly.

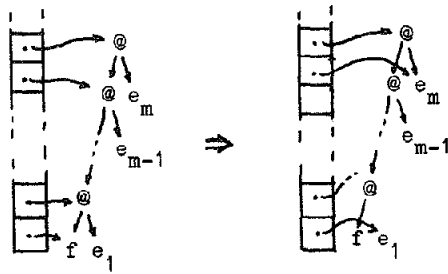


Figure 5. Rearrangement of the stack after unwind.

After the stack rearrangement the function code is executed; see also compilation rule F. If there were too few carried arguments in the application then a premature return is performed, rule 9.

The RET instruction performs a return from EVAL if the function code has updated the apply node for the return value with an integer, boolean, nil or cons node, rule 10. If the updated node is an apply node or function node then the UNWIND state is reentered, to continue the reduction of the new graph; an example when this happens is shown in figure 6 which illustrates reduction of the expression $f(g.\text{nil})\ 3$, where $f\ 1 = \text{hd}\ 1$ and $g\ x = 2*x$. The value of $f(g.\text{nil})$ is the function g , and f has one 'extra' argument supplied. After EVAL and two unwind transitions we have the configuration shown in 6(b), the top of the stack is then made to point to the argument of f , figure 6(c). The code for f then computes the value of $\text{hd}\ 1$, which is the function g , and updates the apply node of the application $f(g.\text{nil})$ with the function node g , figure 6(d). Since the entire graph for which EVAL was called for is not yet fully reduced, the RET 1

instruction of the code for f makes the machine reenter the unwind state, figure 6(e), and the top of the stack is made to point to the argument of g , figure 6(f). The code for g then computes the value of $2*x$ and updates the top apply node with the integer 6. The RET 1 instruction of the function g finally performs a proper return from EVAL, figure 6(h).

The fact that 'extra' carried arguments can be applied to function in this manner, and in general we cannot know in advance how many extra, is the reason for accessing parameters and variables relative to the top of the stack (instead of relative to the bottom which perhaps at first sight would seem more natural).

5.6 Let and letrec expressions

The code for a let or letrec expression constructs the graphs for the locally defined expressions and puts pointers to these graphs onto the stack. When leaving the code for the let or letrec expression these stack elements are removed by the SLIDE instruction; see compilation rules E16, E17 etc. The recursive local definitions in letrec expressions are implemented by constructing cyclic graphs, see scheme Cletrec in table 3. As an example consider the code sequence

```
C[letrec x = f x in x x] r n =
  Cletrec[x = f x] r[x=n+1] (n+1);
  C[x x] r[x=n+1] (n+1); SLIDE 1 =
  ALLOC 1; PUSHFUN f; PUSH 1; MKAP; UPDATE 1;
  PUSH 0; PUSH 1; MKAP; SLIDE 1.
```

Figure 7 shows some of the intermediate machine states when executing this code sequence. To construct the graph for x we must have a pointer to x , for this purpose a HOLE node is allocated by the ALLOC 1 instruction; when $f\ x$ has been constructed the HOLE node is updated with this graph.

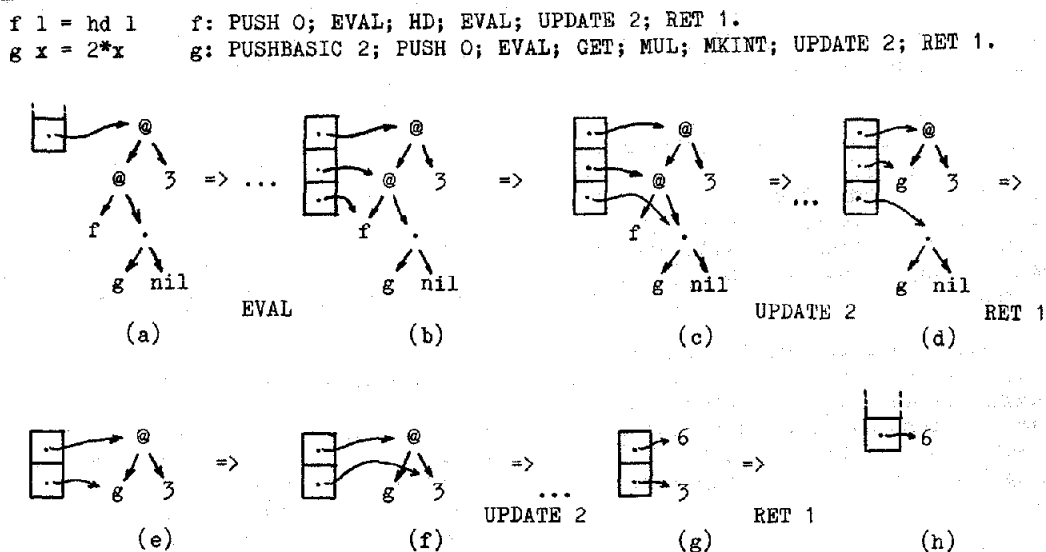
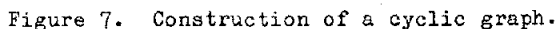


Figure 6. Graph reduction when a function returns a function.



```

B[mul x x] [x=1] 1 =
B[x] [x=1] 1; B[x] [x=1] 1; MUL =
PUSH 0; EVAL; GET; PUSH 0; EVAL; GET; MUL.

```

Here the second eval instruction is clearly useless and can be eliminated. Apart from having removed a useless EVAL instruction, conditions also become better for target code generation from the G-code, since we get longer code sequences unbroken by calls to EVAL and may thus keep things in machine registers a bit longer.

Because of the cost involved in construction and reduction of expression graphs, it is cheaper to evaluate some expressions directly than to construct their graphs, even if the value is not going to be used. This is the case for expressions involving constants, variables which has been evaluated previously, and arithmetic and logical primitive functions (we ignore the problem of overflow and other exceptions). As an example, consider construction of the expression $2*x+y$. The compilation rules in table 3 gives us

```

C [add(mul 2 x) y] [x=2,y=1] 2 =
PUSHFUN add;
PUSHFUN mul; PUSHINT 2; MKAP; PUSH 2; MKAP;
PUSH 3; MKAP.

```

If the variable x has been previously evaluated, it is safe to compute the value of $2*x$, and instead we can use the code sequence

```

PUSHFUN add;
PUSHBASIC 2; PUSH 1; GET; MUL; MKINT;
PUSH 3; MKAP.

```

and if both x and y has been previously evaluated, we can use the code sequence

```

PUSHBASIC 2; PUSH 0; GET; MUL;
PUSH 1; GET; ADD; MKINT.

```

When dealing with expressions with list values the situation is similar. For instance, consider construction of the expression `tl l`, as in the function definition

```
f l = if null l then ... else g (tl l)
```

Because of the test in the condition part of the if-expression, not only can we know for sure that l has been evaluated, in the else part of the if-expression we can also assert that l is on cons form. To construct the expression `tl l`, instead of using

```
PUSHFUN tl; PUSH 2; MKAP
```

we can use the code sequence

```
PUSH 1; TL.
```

Not only does this avoid allocation of an apply node, it also removes the overhead of executing the code for the `tl` function when function `g` calls for evaluation of its argument.

When a variable with a list value cannot be determined statically to be on cons-form, we can test for this dynamically, with instructions MKHD and

MKTL, used in the following compilation rules.

```

C [hd e] r n = C[e] r n; MKHD
C [tl e] r n = C[e] r n; MKTL

```

MKHD and MKTL test whether the top of stack is on cons-form, and if this is the case then behaves as the HD and TL instructions respectively, otherwise constructs the graphs. These instruction are defined by

```

<o, MKHD.c, n.s, v, G[n=CONS n1 n2], E, D> =>
<o, c, n1.s, v, G[n = CONS n1 n2], E, D>
otherwise:
<o, MKHD.c, n.s, v, G, E, D> =>
<o, c, n1.s, v, G[n1=AP n2 n2=FUN hd], E, D>

```

and similarly for MKTL.

The analysis shown above can detect call-by-name to call-by-value transformations only locally within a function. A more general method would be to use a global analysis method, as described in [Mycr80]. A future version of our compiler may include such an analysis phase.

7. Implementation

This section discusses some features of our compiler implementation of the G-machine concept. The source language is a completely function variant of ML [Gord79,Miln84], with call-by-name semantics. The last phase of the compiler translates the G-machine code into target code for the VAX-11 computer.

7.1 Compiler organisation

The compiler is organised into the following parts:

Syntax analysis:

Builds an abstract syntax tree of the program.

Type checking:

Checks that the program is well-typed, using a polymorphic type checking algorithm [Miln78].

Program transformation:

Transforms the program into a set of functions, possibly mutually recursive, as described in section 5.1. Also, the user defined data types and pattern matching is transformed into simpler constructs.

Value analysis:

Performs the analysis on evaluated variables as discussed in section 6.2.

G-code generation:

Translates the functions into G-machine code.

Target code generation:

Translates the G-machine code into assembly code for the VAX-11 computer.

The entire compiler, except for the syntax analysis, has been written in `fc` [Augu82], a functional language with lazy evaluation, a forerunner to the present implementation based on our earlier ideas of compiled graph reduction [John81]. We are currently in the process of rewriting the compiler into its own language.

7.2 Target code generation

For target code generation, the components of the G-machine state is mapped onto the target computer in the following way:

- O is printed on standard output.
- C is the target code of the currently executed function, and the program counter.
- S is a data area for the pointer stack, and a stack pointer register (called ep).
- V is the system stack and stack pointer (sp).
- G is a large heap area divided into two equally sized halves, and a register (called hp) as heap pointer, pointing to the next free location (see below).
- E is the target code for the functions, with code that performs nr-of-arguments-check.
- D is the system stack and stack pointer (sp). Only pointers into S stack and into the system stack are pushed, not entire stacks and dumps as description of the abstract machine suggests.

Both the V stack and the dump D is mapped onto the same stack in the target machine, which is possible because things pushed onto the V stack are only used locally in functions which pushed the value.

The garbage collector is a variant of Fenichel-Yochelson's copying garbage collector [Feni69], but for vary-sized cells, and works as follows. The heap is divided into two equally sized areas. Memory is allocated from one area at a time by simply incrementing the heap pointer hp, and when running out of memory in one area the entire graph is copied into the other heap area, leaving the garbage behind, also updating the pointers on the pointer stack S. In the target code, before an instruction sequence that allocates a certain amount of memory, a check is made if that amount of memory is available on the heap, if not the garbage collector is invoked. A disadvantage of this method of memory management is that only half of the total available memory can be utilised; however on computers with large virtual address spaces this is not a serious problem. To its advantage, the time used for garbage collection is proportional to the size of the graph, (not the size of the heap area, as it is for mark-scan methods) thus taking little time for small graphs.

The target code generation is done by deferring some operations on the pointer stack S and basic value stack V, and instead simulate the contents of the topmost elements. Thus instructions PUSHINT,

PUSHFUN, PUSHBASIC, etc, which pushes constants, will in the code generator push these constants on the simulated stacks. The instruction MKAP, for instance, will thus take two arguments from the simulated stack if nonempty, otherwise from the real stack. To bring out the main idea, a simple example of target code generation is shown in figure 9, which constructs the graph for the expression '3.f 5'. In the simulated stack 'fun f' refers to a pointer to a function node f, 'int i' refers to an integer node with value i, and 'heap n' refers to a pointers into the heap at location n. In the code, newly created nodes on the heap are referred to relative to the hp register, and since node allocation changes the value of hp, we also need to carry along a current relative value of hp, called HP. Function nodes, integer nodes, boolean nodes and the nil node are not allocated each time on the heap; instead pointers to nodes in a constant area are used. (The simulated V stack is irrelevant for this example and is not shown.)

A further possibility which is not shown in this example is to allocate machine registers for entries into the simulated stacks, particularly the V stack for the result of the usual arithmetic operations.

The target code is assembled in the usual manner, and loaded together with the runtime system to make an executable file. The runtime system contains code for PRINT, EVAL, unwind, the garbage collector, and also target code for the primitive predefined functions add, sub etc.

7.3 Performance

We have compared our implementation with a couple of other implementations of functional languages that have been available to us, both with strict and lazy evaluation. The implementations in the table below are the following:

1. Our implementation; lazy evaluation, executes VAX-11 code.
2. Cardelli's ML system [Card84]; strict evaluation, executes VAX-11 code.
3. The Liszt Lisp compiler under UNIX; strict evaluation, executes VAX-11 code.
4. The ML implementation in the LCF system; strict evaluation, interprets Lisp.
5. SASL, based on the SECD machine [Turn75]; lazy evaluation, interpretative.
6. C compiler under UNIX (applies only to the Fibonacci program).

G-code	VAX assembler code	HP	Simulated S stack	Remark
		0	()	Start configuration
PUSHINT 3		0	int 3.()	Push pointer to integer constant 3
PUSHFUN f		0	fun f.int 3.()	Push pointer to function node for f
PUSHINT 5		0	int 5.fun f.int 3.()	Push pointer to integer constant 5
MKAP	movl \$APPLY,(hp)+	4		Tag of apply node to heap ...
	movl \$C F,(hp)+	8		Fun. part = fun f to heap ...
	movl \$I 5,(hp)+	12	heap 0.int 3.()	Arg. part = int 5 to heap.
CONS	movl \$CONS,(hp)+	16		Tag of cons node to heap ...
	movl \$I 3,(hp)+	20		Head part = int 3 to heap ...
	moval -20(hp),(hp)+	24	heap 12.()	Tail part = result of MKAP to heap ...
	moval -12(hp),-(ep)	24	()	Move result to real S stack.

Figure 9. Target code generation from graph construction code.

The table below shows the execution time in seconds for three programs: fib(20) using fib(n)=if n<2 then 1 else fib(n-1)+fib(n-2), primes up to 300 using sieve of Erathosthenes, and insertion sort of 100 random elements.

	1.	2.	3.	4.	5.	6.
Fibonacci	0.92	0.5	1.1	46	31	0.46
Primes	0.5	1.2	1.1	29	20	-
Insert sort	0.37	1.0	0.8	15	12	-

The programs above have been chosen so that the results are the same independent of whether lazy or strict evaluation is used, but in general lazy evaluation permits a more direct programming style. It should be noted that in our Fibonacci program, in the recursive call to fib the arguments are passed by value, due to the analysis on evaluated variables described in section 6.2.

8. Related work

Jones and Muchnick [Jones82] gives an alternative evaluation mechanism for combinator expressions, with a compilation algorithm which translates combinators to fixed-program code for a stack machine.

Hudak's combinator based compiler [Hudak84] resembles our work in many respects. He uses the standard combinators as a convenient intermediate language for performing program transformations and optimisations. The program is the converted into one containing macro-combinators, which is similar to Hughes' super-combinators [Hugh82] and our global function definitions. Each macro-combinator is then translated into code for a conventional machine.

Dick Kieburtz et. al at Oregon Graduate Center is currently in the process of designing and implementing a VLSI chip for the G-machine.

9. Acknowledgements

This work was supported by the Swedish Board for Technical development (STU). The compiler has been implemented together with Lennart Augustsson, and many ideas has grown out of this close cooperation. I also wish to thank Alan Mycroft for helpful comments on earlier drafts of this paper, and the members of the Programming Methodology Group for numerous nice cake parties.

10. References

- [Augu82] L. Augustsson, "FC manual", Memo 13, Programming Methodology Group, Chalmers University of Technology, Goteborg (1982).
- [Back78] J. Backus, "Can Programming be Liberated from the von Neumann Style? A functional style and its algebra of programs", *Comm. ACM*, Vol. 21 no. 8 pp. 613-641 (August 1978).
- [Card84] L. Cardelli, "ML under UNIX", *Polymorphism: The ML/LCF/Hope Newsletter*, Vol. 1 no. 3 (January 1984).
- [Feni69] R. Fenichel and J. Yochelson, "A LISP garbage-collector for virtual memory computer systems", *Comm. of the ACM*, Vol. 12 no. 11 pp. 611-612 (November 1969).
- [Frie76] D. P. Friedman and D. S. Wise, "Cons should not evaluate its arguments", pp. 257-284 in *Automata, Languages and Programming*, Edinburgh Univ. Press (1976).
- [Gord79] M. Gordon, R. Milner, and C. Wadsworth, "Edinburgh LCF", *Lecture Notes in Computer Science*, Vol. 78, Springer Verlag (1979).
- [Huda84] P. Hudak and D. Kranz, "A Combinator-based Compiler for a Functional Language", pp. 122-132 in *Proc. 11th ACM Symp. on Principles of Programming Languages* (1984).
- [Hugh82] J. Hughes, "Super Combinators: A New Implementation Method for Applicative Languages", pp. 1-10 in *Proceedings of the 1982 ACM Symposium on Lisp and Functional Programming*, Pittsburgh (1982).
- [John81] T. Johnsson, "Code Generation for Lazy Evaluation", Memo 22, Programming Methodology Group, Chalmers University of Technology, Goteborg (1981).
- [Jones82] N.D. Jones and S.S. Muchnick, "A Fixed-Program Machine for Combinator Expression Evaluation", pp. 11-20 in *Proceedings of the 1982 ACM Symposium on Lisp and Functional Programming*, Pittsburgh (1982).
- [Land64] P. J. Landin, "The Mechanical Evaluation of Expressions", *Computer Journal*, No. 6 pp. 308-320 (January 1964).
- [Land66] P. J. Landin, "The Next 700 Programming Languages", *Comm. of the ACM*, Vol. 9 no. 3 pp. 157-164 (March 1966).
- [Miln78] R. Milner, "A Theory of Type Polymorphism in Programming", *Journal of Computer and System Sciences*, Vol. 17 no. 3 pp. 348-375 (1978).
- [Miln84] R. Milner, "Standard ML Proposal", *Polymorphism: The ML/LCF/Hope Newsletter*, Vol. 1 no. 3 (January 1984).
- [Mycr80] A. Mycroft, "The Theory and Practice of Transforming Call-by-Need into Call-by-Value", pp. 269-281 in *Proc. 4th Int. Symp. on Programming, Lecture Notes in Computer Science*, Vol. 83, Springer Verlag, Paris (April 1980).
- [Turn75] D. A. Turner, "An implementation of SASL", TR/75/4, St. Andrews (1975).
- [Turn79] D. A. Turner, "New Implementation Techniques for Applicative Languages", *Software - Practice and Experience*, Vol. 9 pp. 31-49 (1979).