

The Structure of the Lisp Compiler

In the course of my work in algebraic symbol manipulation, I have had the pleasure to implement a Lisp system with compiler for several S/360 operating systems. My confederates in this venture were J. Griesmer, J. Harry, and M. Pivovonsky. In the course of that development I had the opportunity to study the history and structure of the Lisp compiler. It is my intention to pass along some of that information.

The first Lisp compiler was written by Robert Brayton with the assistance of David Park, in SAP for the 704. That compiler was started in 1957 and was working in 1960 by which time Brayton left MIT. During that interval of time a Lisp compiler written in Lisp was implemented by Klim Maling but that compiler was apparently dropped. The argument advanced was that Brayton's being written in assembly language, would obviously be faster. Difficulties in maintenance developed when Brayton left the project. After Brayton and Maling, Timothy Hart and Michael Levin wrote a compiler in Lisp which was distributed with the 704 Lisp 1.5 system. The compiler that I am most familiar with and will describe today is a descendant of that compiler.

The Lisp compiler has developed over the years due to

several factors:

1. The Lisp language has been refined and the compiler has had to keep up.
2. The Lisp system was implemented on new computers with different order codes and memory organizations, therefore the code generation had to change.
3. Bugs were identified and corrected.
4. Various optimizations of the compiler itself and the code it produces were deemed desirable.

To date none of these factors have completely receded into the woodwork and therefore development continues.

To cite just a few of the more sophisticated and active efforts which I consider lambda calculus motivated:

GEDANKEN by J. Reynolds; PAL by A. Evans; CPL - Barron et. al.; ECL by B. Wegbreit.

To these should be added the many straight Lisp 1.5 type efforts, of which the Stanford Lisp systems, the BBN-Lisp on PDP-10 with the TENEX time sharing system, BALM on the CDC 6600 by Harrison, et al, are just a few examples. S/360 Lisp is of this latter type and while it is likely to persist for some years could hardly be touted as a penultimate system. I foresee a series of modest improvements. The overriding consideration being improvement but compatibility with the

large body of code that we currently have available.

It is with all this in mind that I launch into this presentation of our S/360 Lisp system and the structure of its compiler.

Firstly, the Lisp system should be understood as having:

1. A language
2. A data-model including storage organization and management
3. A supervisor
4. A large set of built in primitives
5. Two evaluation mechanisms, the interpreter and the compiler
6. A definition schemata

All these facets are of course interesting and vital. But today I am concerned only with the compiler and will to some extent, assume that you are familiar with the other facets of the Lisp system. Moreover that you are here mainly to learn about the organization of the compiler and are unconcerned about other topics which I would be otherwise happy to discuss.

The Lisp compiler is a run time function which creates named code bodies. The effect is that new primitive operators are introduced into the system. Such operators may have access to global variables and are free to have global

side-effects. The compiler is written in Lisp. It is rather well understood and indeed correctness proofs of several subset compilers have been published in 1971 by R. London.

As we are concerned with an extremely simple language with typeless data objects, the compiler is not overly complex. The compiler is primarily concerned with variable binding and variable evaluation, and to that end it uses more efficient environment models than the a-list of the interpreter. These environment models are in order of efficiency: stack variables, SPECIAL variables and COMMON variables.

It is a goal of the definition schema in this lisp system that functions may each be defined in any order. This is a real boon to the interactive definition process, but could lead to a class of errors or else result in inefficient code. The errors that I have in mind come when the number of arguments or their evaluation or transmission modes is changed in a redefinition.

In fact that is just what happens and furthermore is considered tolerable. For reasons of efficiency the compiler prefers to think that all bound variables are purely local (unless otherwise informed) for which it has an efficient stack model. As you can see this is an incorrect assumption which allows for more efficient code but makes the programmer identify those variables which may be used free by some called function.

not clear to me

Another incorrect assumption that this compiler thrives on is that a call to a currently undefined function is a call to a function that expects exactly the number of arguments given and that it expects them evaluated. This is not so bad in the case the function remains undefined as a run time error occurs. Also any attempt to define it as other than the assumed class can give rise to an error.

The compiler thus assumes that the user declares his free variables correctly and compiles his unusual functions before any calls to it are compiled. The three modes for variables that the compiler uses are:

1. stack variables.
2. SPECIAL variables.
3. COMMON variables.

The implementations of each of these are:

1. Stack variable -- A stack cell is reserved for each bound variable occurring in the lambda-expression, these cells are filled in at entry to the code body with the arguments of the function. A stack cell is allocated for each program variable found as bound variables in embedded PROGS, these are initialized to NIL at function entry. Temporary variables are created by the compiler itself. All of these cells are accessible only to the current activation of the code body in question.

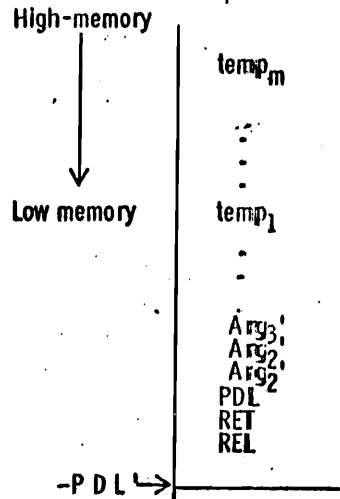


Fig. 1

2. SPECIAL variables -- This mode is used for those variables declared SPECIAL at compile time. A value cell is associated with the identifier which denotes a variable. At binding time the contents of the stack cell allocated to that bound variable is exchanged with the contents of the value cell. Subsequent references and assignments are to the value cell. At function return time these two cells are exchanged again, thus restoring the former values.

The implementation of SPECIAL variables is quite efficient but inadequate in the following cases.

1. The downward FUNARG:

A functional argument is passed which uses a free variable which has been bound again in the environment the functional value is applied.

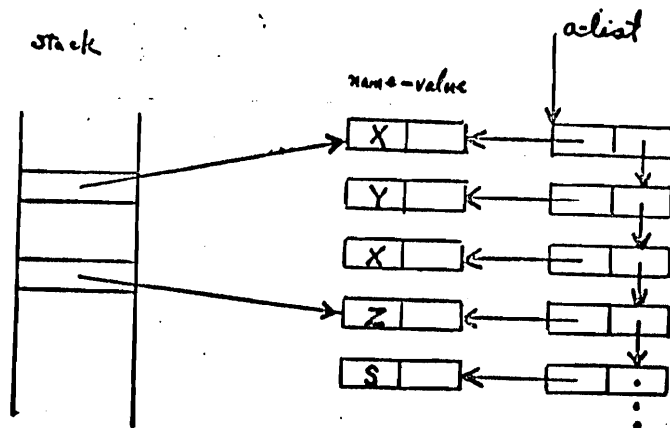
2. The upward FUNARG:

A functional value is returned which has a free variable which has a different binding when that value is applied.

For these cases the full power of the COMMON variable mode is necessary. The Lisp interpreter has access to SPECIAL values only if no variable by that name occurs on the a-list.

3. COMMON variables -- This is the most correct variable model used because it uses the a-list in a manner isomorphic to the use by EVAL. The embodiment is to have a stack cell point to the name-value pair on the a-list. Evaluation of such a variable is a LOAD and a CDR. Assignment is a LOAD and a RPLACD. When COMMON variables are bound (activated), the a-list is augmented by appending that name-value relationship cell. At deactivation time the a-list is restored to its former value. In those functions which use a COMMON variable free, or for that matter any undeclared operand free variable, the a-list is searched at binding (activation) time and a stack cell is made to point to the appropriate name-value cell.

2/67
?



Compiler treatment of common variables Z and X.

Fig.

The COMMON variable mechanisms of the compiled code are more efficient than those of the interpreter except that free variables which are mentioned but not used do get extra attention. The COMMON mechanism is less efficient than the stack or SPECIAL variable mechanisms. It is more correct however.

The compiler treats ordinary operator names like free variables and compiles the code to:

1. Evaluating the arguments.
2. Transmit the arguments.
3. Do a subroutine jump indirectly through the SPECIAL value cell of that identifier.

It is the users responsibility to see that any operator that deviates from the ordinary class receive properties which the compiler may know before any function which calls it is compiled.

Many of these defaults and responsibilities were put in as aids to the interactive program development, but are in fact the sources of some errors.

It is interesting to note that in Brayton's early compiler, the "compile before or along with" conventions were rigidly enforced. A function not satisfying those requirements

didn't compile.

Despite my earlier promise that I would concentrate strictly on the structure of the compiler I should like to describe the Lisp language. My argument is that even those of you who are expert at Lisp need to know precisely what Lisp I am concerned with.

To that end I present a syntax description of Lisp as I see it.

Syntax of LISP

e is: $\{c \mid \text{id}\}$ a constant or variable name

$\left\{ \begin{array}{l} \text{LAMBDA} \\ \text{FLAMBDA} \end{array} \right\} \left\{ \begin{array}{l} (\text{id}^*) \\ \text{id} \end{array} \right\} \text{body}$. where body is an e

(LABEL id e)

(SETQ id e)

(COND (p q)*) where p, q are e

(PROG (id*) s*)

where s is

$\left\{ \begin{array}{l} \text{Id} \\ (\text{GO id}) \\ (\text{RETURN e}) \\ (\text{COND (p s)*}) \\ e \end{array} \right\}$

(FUNCTION e)

(QUOTE s-exp)

(rator rand*)

where rator, rand are e.

A Lisp expression e is:

c a constant.

id a variable name or identifier.

(lambda-op bv body) a lambda-expression or procedure
 where lambda-op is LAMBDA or FLAMBDA, and
bv the bound variables is either id or (id*), and
body is an e.

(LABEL id e) a label-expression.

(SETQ id e) assignment.

(COND (p g)* a conditional-expression
 where the predicate p is an e, and
 the consequent g is an e.

(PROG pv s*) a program-expression
 where pv the program-variable part is (id*), and
 each statements s is a:

label which is an id, or

(GO id) a go-statement, or

(RETURN e) a return-statement, or

an e, or

(COND (p s)* a conditional-statement.

(FUNCTION e) .

(QUOTE s-expression) a quoted s-expression.

(rator rand*) a combination
 where the operator rator is an e, and
 each operand rand is an e.

It should be noted that except for constants and variables

every Lisp expression is a combination. Some of these combinations are distinguished for semantic reasons.

Semantics in Lisp is usually given by describing an interpreter (EVAL e a) where e has been described and a is an environment which gives meaning to free variables. The environment a is an ordered set of name-value relationships which is most generally represented as an a-list.

As an aside I should like to mention my own bias that an interpreter described as a S,E,C,D transform is even more illuminating. Such transforms were described by P. Landin in 1964. (see appendix A for such a semantics presentation)

It is well known that Lisp requires a minimal set of data primitives for its own self-description. To define EVAL in Lisp we require only the following additional primitives:

(EQ x y) identity.

(ATOM x) the atom predicate.

(IDENTP x) the identifier predicate.

(CAR x) and (CDR x) the selectors.

(RPLACD x y) and (RPLACA x y) the storing functions.

and an allocation function like CONS.

The fact that very little is required to understand Lisp and the fact that it is self described and self implemented has allowed the users to play with its semantics. Lisp has not been very self protective. Some examples of such introduced semantics are macro-expressions and code-expressions.

A macro-expression m-e is simply a combination whose rator is understood to be the name of a macro definition function mdef. The value of a macro-expression is given by: (EVAL (mdef m-e) a) .

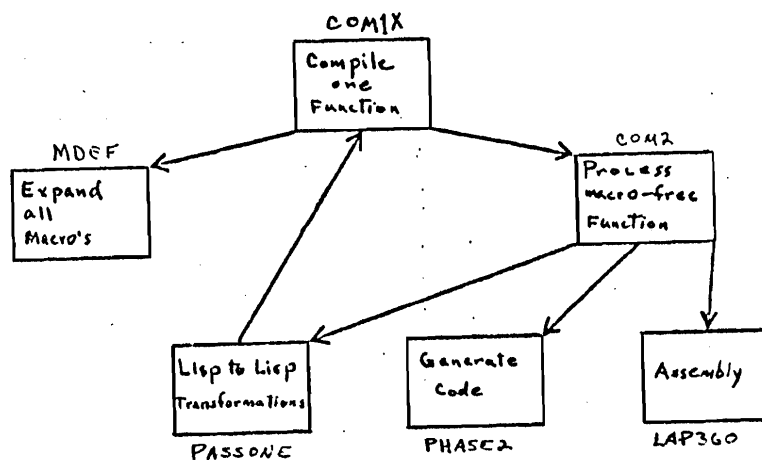
A code-expression is a combination of the form:

(CODE instr*) where the instr are LAP instructions.

Both macro and code expressions are intended as compilation tricks but EVAL could handle them. In practice EVAL handles macro's but not code-expressions.

Skipping any further discussion of semantics I believe we should be in a position to follow a structural description of the compiler.

Each named procedure is compiled by a Lisp function COM1X which has the following structure:

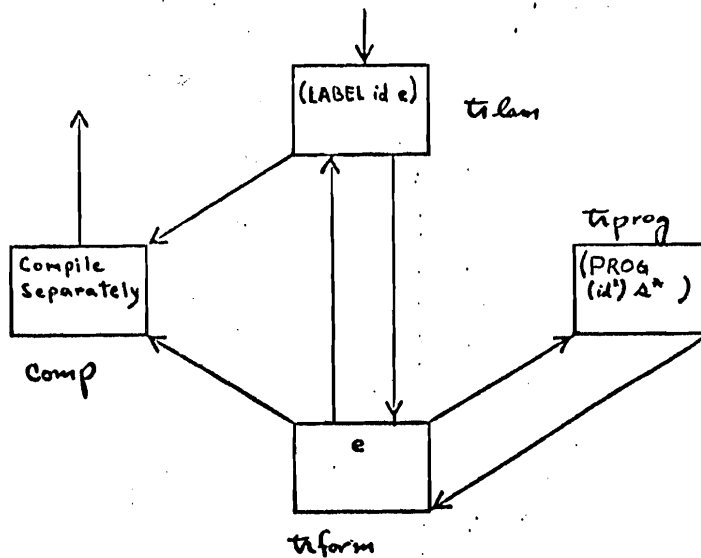


COM1X calls MDEF which expands all the Lisp macros. COM1X is recursively called for certain procedure-valued expressions. In which case a separate function is created and then referenced by the "containing" function.

COM2 controls the compilation of a properly macro-expanded named lambda-expression. There are three parts to this compilation namely, PASSONE, PHASE2, and LAP360.

PASSONE performs such Lisp to Lisp transformations as:

1. Recursion removal
2. Free variable identification
3. Introduction of binding mechanisms
4. Elimination of indefinite number of arguments
5. Elimination of FLAMBDA operators
6. Remove operand procedure-valued expressions
7. Introduce interpreter calls for delayed evaluation conditions



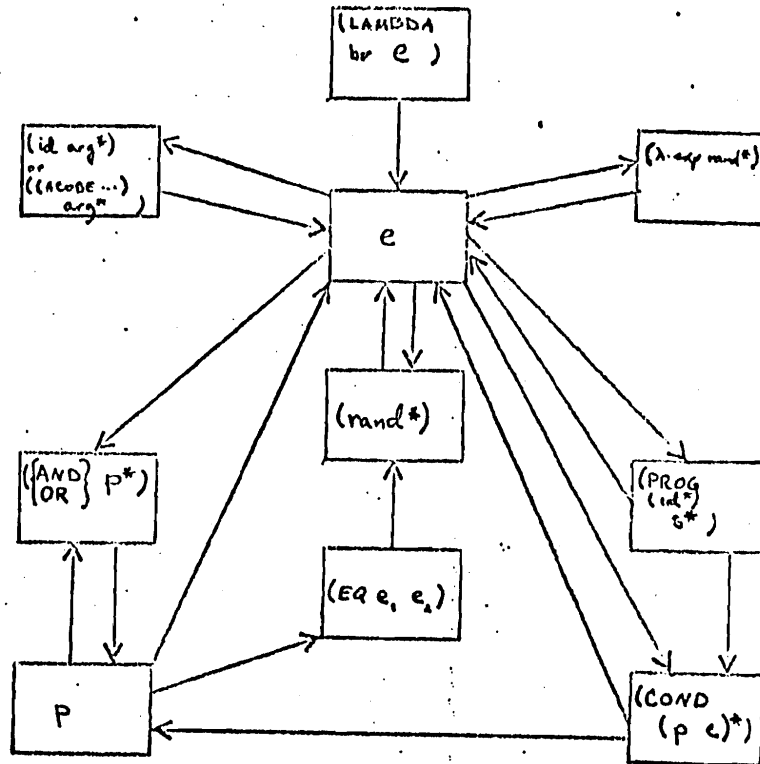
PASSONE

Fig (see appendix B for greater detail on PASSONE.)

PHASE2 receives a much simplified lambda-expression from PASSONE. It produces a list of LAP (List Assembly Program) instructions which is meant to be the execution equivalent of the lambda-expression. PHASE2 creates the code. It knows what to do for:

(QUOTE s-expression)	constants
(SPECIAL id)	
(SETQ id e)	assignment
(COND (p q)*)	conditionals
(PROG pv s*)	program blocks
(RETURN e)	program returns
(GO id)	go-tos
(AND e*)	
(OR e*)	
(*CODE instr*)	code-expressions
((LAMBDA (id*) e) e)	lambda-exp operator
(id e*)	calls

PHASE2 also generates the code for variable binding prologues and variable restores on value return epilogues.



(see appendix C for greater detail on PHASE2)

LAP360 is a very simple minded 2-pass assembler. The only interesting fact worth mentioning is that it finally resolves special variables, quotes and function calls.

References

Reynolds, J. C., GEDANKEN, A Simple Typeless Language Based on the Principle of Completeness and the Reference Concept. CACM 13, 5 (1970) 308-319

Evans, A., PAL - a language designed for teaching programming linguistics, ACM 23rd National Conference, August 1968.

Barron, D., et al., The main features of CPL, Computer Journal 6, 2, July 1963.

Landin, P., The Mechanical Evaluation of Expressions, Computer Journal 6, 4, 1964.

London, R. L., Correctness of two compilers for a Lisp subset, Artificial Intelligence Memo AIM-151, Stanford University, October 1971.

Wegbreit, B., The ECL programming system, Aiken Computation Laboratory (Harvard University, 1971 (in preparation)).

APPENDIX A. Semantics ---

965 - IBM - 04

965 - IBM - 04

965 - IBM - 04

965 - IBM - 04

APPENDIX B. PASSONE transformations ---

965-IBM-04

965-IBM-04

Iterative form recursion removal

```

triam [ (LABEL x (LAMBDA (id1 ... idn)
  (COND { (p q) | (pi (x ai1 ... ain)) } *)); b ]
  (P**G (g1 ... gn)
  gs (COND { (p (RETURN q)) | (pi (GO gi)) } *)
  { gi (SETQ gi ai1) ... (SETQ gn ain) (GO gt) }
  gt (SETQ id1 g1) ... (SETQ idn gn) (GO gs))); b ]
  
```

Constant forms

$trform [var ; b]$ where var is on id list and not declared special.
 $= (COMMON var)$ and var is added to the last free.
 and
 $free := (var \cdot free)$

if var were declared special then
 $= (SPECIAL var)$

if $var \in b$ but declared:

$\left\{ \begin{array}{l} (COMMON var) \\ (SPECIAL var) \end{array} \right\}$ if common
 if special

if var not declared but $var \in b$
then var

```

triam [ (LABEL x (LAMBDA (bv (PROG ...))); b ]
  (LABEL x (LAMBDA (bv (P**G ...))); b ]
  
```

```

triam [ (LABEL x (LAMBDA (bv body)); b ]
  (LABEL x (LAMBDA (bv { trform [ body; bv: b ] /
  (PROG pv
    alist-save
    common-var-bind
    free-var-bind
    special-var-bind
    (SETQ gv trform [ body; bv: b ] )
    special-restore
    alist-restore
    (RETURN gv) ) } ) )
  
```

where $pv = (\{ g_a | empty \} g_v (free^*)$

$trform [T ; b] = (QUOTE *T*)$
 $trform [F ; b] = (QUOTE NIL)$
 $trform [x ; b] = (QUOTE x)$ if $atom[x] \wedge \sim idlist[b]$
 $trform [(QUOTE x) ; b] = (QUOTE x)$

$\text{tiform}[(\text{id} \text{ rand}^*); b]$

where id is undefined and undeclared and free,
or is defined as \Rightarrow SUBR or EXPR.

$= (\text{id} \text{ tiform}[\text{rand}; b]^*)$

$\text{tiform}[(\text{id} \text{ rand}^*); b]$

where id is defined as FSUBR or FEXPR

$= (\text{id} (\text{QUOTE} \text{rand})^*)$

$\text{tiform}_1[(\text{id} \text{ rand}^*); b]$

where id is defined as SUBR₁ or EXPR₁

$= (\text{id} (\text{LIST} \Rightarrow \text{tiform}_1[\text{rand}; b]^*))$

$\text{tiform}[(\text{id} \text{ rand}^*); b]$

where id is defined as FSUBR₁ or FEXPR₁

$= (\text{id} (\text{LIST} (\text{QUOTE} \text{rand})^*))$

$\text{tiform}[(\text{id} \text{ rand}^*); b]$

where id is either $\in b$ or special or common.

$= (\text{APPLX} \text{tiform}[\text{id}; b] (\text{LIST} \text{tiform}[\text{rand}]^*))$

note: assumes arguments evaluated.

$\text{trform}[(\text{SETQ} \text{id} \text{e}); b]$

$= \left\{ \begin{array}{l} (\text{RPLACD} \text{id} \text{e}') / (\text{SETQ} (\text{SPECIAL} \text{id}) \text{e}') / \\ (\text{SETQ} \text{id} \text{e}') \end{array} \right\}$ where $\text{e}' = \text{trform}[\text{e}; b]$

$\text{trform}[(\text{COND} (\text{p} \text{q}^*); b]$

$= (\text{COND} (\text{trform}[\text{p}; b] \text{trform}[\text{q}; b]^*))$

$\text{trform}[(\text{GO} \text{id}); b] = (\text{GO} \text{id})$

$\text{trform}[(\text{AND} \text{p}^*); b] = (\text{AND} \text{trform}[\text{p}; b]^*)$

$\text{trform}[(\text{OR} \text{p}^*); b] = (\text{OR} \text{trform}[\text{p}; b]^*)$

$\text{trform}[(\text{FUNCTION} \text{e}); b] = (\text{FUNC} (\text{QUOTE} \text{comp}[\text{e}; \text{g}]) \$\text{ALIST})$

$\text{trform}[(\text{*CODE} \text{instr}^*); b] = (\text{*CODE} \text{instr}^*)$

$\text{trform} [(P^{**}G \text{ pv } s^*) ; b]$
 = $\text{trprog} [(PROG \text{ pv } s^*) ; b]$

 $\text{trform} [(PROG (id_1 \dots id_n) s^*) ; b]$
 = $\text{trprog} [(PROG (id_1 \dots id_n) (SETQ \text{ id}, (QUOTE \text{ NIL}))$
 $\dots (SETQ \text{ id}_n (QUOTE \text{ NIL})) s^*) ; b]$

 $\text{trform} [\lambda\text{-exp} ; b] = (\text{FUNC } (QUOTE \text{ comp} [\lambda\text{-exp} ; g]) \$ALIST)$

$\text{trprog} [(PROG \text{ pv } s^*) ; b]$

= (PROG pv'
 alist-save
 common-var-bind
 free-var-bind
 special-var-bind
 (SETQ g_v (PROG $\text{pv}' s'$))
 special-restore
 alist-restore
 (RETURN g_v))

where $\text{pv}' = (g_v \text{ coms}^* \text{ spec}^* \text{ free}^*)$
 $\text{pv}' = \text{pv} - [\text{coms}^* \cup \text{spec}^* \cup \text{free}^*]$

$\text{trform} [((LABEL \text{ id } e) \text{ rand}^*) ; b]$
 = $\text{trform} [(\text{trlam} [(LABEL \text{ id } e) ; b] \text{ rand}^*) ; b]$

$\text{trform} [((FUNCTION e) \text{ rand}^*) ; b]$
 = $\text{trform} [(e \text{ rand}^*) ; b]$

$\text{trform} [((QUOTE s\text{-exp}) \text{ rand}^*) ; b]$
 = $\text{trform} [(s\text{-exp} \text{ rand}^*) ; b]$

$\text{trform} [((*CODE instr^*) \text{ rand}^*) ; b]$
 = $((CODE \text{ instr}^*) \text{trform} [\text{rand} ; b]^*)$

$\text{trform} [((rator \text{ rand}^*) \text{ rand}'^*) ; b]$
 = (APPLX $\text{trform} [(rator \text{ rand}^*) ; b]$
 (LIST $\text{trform} \text{rand}' ; b \cdot$))

$\text{trform} [((LAMBDA (id^*) \text{ body}) \text{ rand}^*) ; b]$
 = $(\text{trlam} [(LABEL g (LAMBDA (id^*) \text{ body})) ; b]$
 $\text{trform} [\text{rand} ; b]^*)$

$\text{trform} [((LAMBDA \text{ id } \text{ body}) \text{ rand}^*) ; b]$
 = $(\text{trlam} [(LABEL g (LAMBDA (\text{ id}) \text{ body})) ; b]$
 (LIST $\text{trform} \text{rand} ; b \cdot$))

$\text{trform} [((FLAMBDA (id^*) \text{ body}) \text{ rand}^*) ; b]$
 = $(\text{trlam} [(LABEL g (LAMBDA (id^*) \text{ body})) ; b]$
 (QUOTE rand^*))

$\text{trform} [((FLAMBDA \text{ id } \text{ body}) \text{ rand}^*) ; b]$
 = $(\text{trlam} [(LABEL g (LAMBDA \text{ id } \text{ body})) ; b]$
 (QUOTE (rand^*)))

APPENDIX C. PHASE2 code generations ---

comval [id; stomap; name]

= (L AC locate [id] [stomap])

where stomap = ((id_j (PDL 4k)) (id_{j-1} (PDL 4[k-1])) ...)

(PDL 4k) is the stack address for the local bound variable id_j.

comval [(QUOTE s-exp); stomap; name]

= (L AC (QUOTE s-exp))

comval [(SPECIAL id); stomap; name]

= (L AC (SPECIAL id))

comval [((LAMBDA (id₁... id_n) body) e₁ ... e_n); stomap; name]

= (comval [e₁; stomap; g₁] :
(ST AC locate [id₁] [stomap']) :

comval [e_n; stomap; g_n] :
(ST AC locate [id_n] [stomap']) :

comval [body; stomap'; name])

where stomap = ((id_j (PDL 4k)) (id_{j-1} (PDL 4[k-1])) ...)

and stomap' =

((id_n (PDL 4[k+n])) ... (id₁ (PDL 4[k+1])) (id_j (PDL 4k)) ...)

comval [(id rand_i*) ; stomap; name] for i > 2

= (comval [rand₁; stomap; g₁] :

(ST AC locate [g₁] [stomap']) :

comval [rand₂; stomap; g₂] :

(ST AC locate [g₂] [stomap']) :

comval [rand_n; stomap; g_n] :

(ST AC locate [g_n] [stomap']) :

(L AC locate [g₁] [stomap']) :

(L MQ locate [g₂] [stomap']) :

(LA PDL 0 locate [g₃] [stomap']) :

(CALL id n))

comval [(COND (p_i e_i)*); stomap; name]

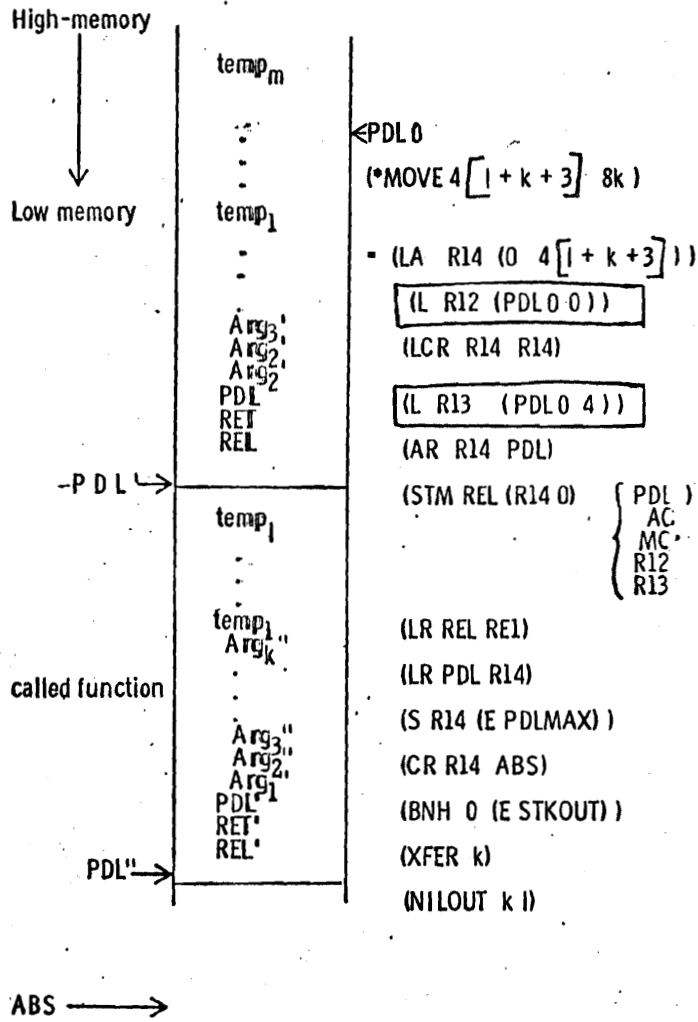
= ({ comval [p_i; stomap; g_i] :
(CR AC NIL) :
({BNE} 0 (LABEL g_{i,1})) :
comval [e_i; stomap; name] :
(B O (LABEL name)) :
g_{i,1}
...
for i = 1 ... *
name)

comval [(PROG pv s*); stomap; name]
 where stomap = ((id_j (PDL 4k)) (id_{j-1} (PDL 4[k-1])) ...)
 = comprog [(s*); pv; name]

- I. (s') is created from (s*) by substitution of g_i for each label_i in (s*)
- II. stomap' is created from pv = (id₁ id₂ ...) and stomap.
- III. The code for each s' in (s') is created:
 - if atom [s'] → s'
 - if s' = (GO label_i) → (B O (LABEL g_i))
 - if s' = (COND (p s*)) → comcond [((ps)*); NIL]
 - else comval [s'; stomap'; g]
- IV. The label name is placed at the end of all this code.

= (c₁ : c₂ : ... c_n name) where c_i is a list of code for s_i.

STACK PLAN 72 - 2



(XFER k)

if $k < 5 \rightarrow ()$ else (MVC $\underline{\quad}$ (PDL 28) $4 \underline{[k-4]}$ (PDL 0, 16))

(NILOUT k 1)

if = 0 $\rightarrow ()$ = 1 \rightarrow (ST NIL (PDL $4k + 12$))= 2 \rightarrow (ST NIL (PDL $4k + 12$)) (ST NIL (PDL $4k + 13$))else \rightarrow (ST NIL (PDL $4k + 12$))(MVC $\underline{\quad}$ (PDL $4k + 12$) $4 \underline{[1-1]}$ (PDL $4k - 16$))