

SISP/I AN INTERACTIVE SYSTEM ABLE TO
SYNTHESIZE FUNCTIONS FROM EXAMPLES

Jean-Pierre JOUANNAUD FRANCE
Maitre-Assistant a l'Institut de Programation
Universite Paris VI
A, Place Jussieu
75005 PARIS

Gerard GUIHO FRANCE
Maitre de Conference - Laboratoire de Recherche
en Informatique
Universite Paris Sud
91405 ORSAY

Jean-Pierre TREUIL FRANCE
Chercheur - Laboratoire de Recherche en
Informatique
Universite Paris Sud
91405 ORSAY

The research presented in this paper is supported
by IRIA-CESORI under contract number 76.

ABSTRACT

SISP/i is an interactive system whose goal is
the automatic inference of LISP functions from a
finite set of examples $\{(x., f(x.))\}$ where $x.$ is a
list belonging to the domain of the function f we
want to infer. SISP/I is able to infer the recur-
sive form of many linear recursive functions and
its stop-condition. SISP/I tries to work with one
example only. When it fails, it asks for new ones:
using then a method of generating new partial sub-
problems, SISP/I is able to perfect its generated
recursive function until it gets a correct one.

I. INTRODUCTION

In this paper we describe the system SISP/I
whose goal is the automatic inference of LISP func-
tions from a finite set of examples $\{(x., l(x.))\}$,
where $x.$ is a list belonging to the domain of the
function f we want to infer.

The problem originates from a more general
one: how to build a "Learning-Question-Answering-
System" (L.Q.A.S.) using a functional method to
provide an answer to any given question. The meth-
od we propose in SISP/I is naturally well adap-
ted to the L.Q.A.S. we are developping (6.1, 17 1.

In the field of "Automatic Programming from
Examples", an important piece of recent work is
THESYS by SUMMERS L5J. The major result of this
work, is the following: using a small number of
well chosen examples
 $\{(NIL, f(NIL)), ((A), f((A)))\dots\}$ THESYS is able
to infer a recursive expression $\$$ equivalent to f
for every x belonging to the domain of f .

Only a small class of functions can however
be obtained by Summers's method, which works by loo-
king for a recurrence relation between representati-
ve predicates $p.$ of the given input structure and
for a recurrence relation between the map functions
 $m.$ providing the given outputs from the given in-
puts. Then, using a fixed point theorem, V is cons-
tructed.

Although Summers's method is very powerful it

has four important drawbacks:

- 1.- The constructed expression $\langle p$ is necessarily
recursive: for instance the identity function will
be infered by $V(x) \gg$ if $X * NIL$ then NIL
else $CONS(CAR(x), ^{(CDR(x))})$
- 2.- THESYS needs well chosen examples, which in
particular must contain the stop condition of the
recursive function V . For instance, the construc-
tion of the function REVERSE requires the following
set of examples:
 $\{(NIL + NIL), ((A) - (A)), ((A B) - (B A)),$
 $((A B C) + (C B A))\}$
- 3.- The function to be constructed has to present
only one "iterative level"¹. For instance, THESIS
fails to construct a correct function corresponding
to the example: $(P Q R S) \rightarrow (P P Q P Q R P Q R S)$.
- 4.- When THESYS has to solve a difficult problem,
it does not try to generate a partial, simpler
problem for which it could either find a correct
solution or perhaps use a knowledge previously
stored in a data base by the system itself. Thus,
THESYS cannot be efficiently used in a L.Q.A.S.
without important modifications.

The method we propose in this paper is very
different in particular, it has the built capacity
to use a Professor in interactive mode. It does
not lie yet on any theoretical groundwork, but allows
us to overcome some of the previous drawbacks, al-
though new ones appear:

- recursion is not automatically infered by the
synthesis algorithm; for instance, using the exam-
ple $((A B C) -* (A B C))$, SISP/I inferes the function
 $\Phi \forall^p(x) = x$ for any x .
- for some "simple" functions, SISP/I needs only
one example $(x, f(x))$.
In the case where a recursive expression is infered,
the stop condition is then found by SISP/I itself.
However the list x must be long enough to be re-
presentative of the function f . For instance,
REVERSE is obtained using the only example
 $((A B C D) + (U C B A))$, but is not obtained with
 $((A B C) > (C B A))$.
- when the function f is "more complicated" SISP/I
fails to construct a correct function with only one
example and it then tries to work with two examples.
- when the function f is "much more complicated",
SISP/I generates a new partial simpler problem
 $(y \gg \&(y))$ where y is defined in terms of x and
 $g(y)$ is defined in terms of $f(x)$. To solve this
new problem, SISP/I sometimes needs a new example
 $(x', f(x'))$ which is used to deduce an example
 $(y' \gg g(y'))$ « The interaction is only used in the
sense of asking for new examples, when necessary.
SISP/I is thus extensible and has the potentiality
to use a self constructed knowledge data base.

Some objections can be raised to our interac-
tive method:

- when a function f needs several examples to be
infered, the professor sometimes has to give an
appropriate sequence of examples.
- we do not exactly know the class of functions
which SISP/I is able to infer. However, it seems
to be much larger than THESYS one. For instance
 $((P Q R S) \rightarrow (P P Q P Q R P Q R S))$ is infered by
SISP/I using only one example whereas the HALF
function $((P Q R S T U) + (P Q R))$, which is infered
by THESYS, requires two examples by SISP/I. In fact,

we hope that SISP will be able to infer a larger class of linear recursive functions.

11. GENERAL DESCRIPTION OF THE METHOD

1.- L\$ngu\$ge

SISP/1 infers functions defined on character strings "ABCD..." which will be represented by the list (A B C D...).

SISP/1 synthesizes LISP-functions built with the following basic functions, described here by examples:

LCAR: (A B C D) -* (A) CDR: (A B C D) -> (B C D)

LRAC: (A B C D) +> (D) RDC: (A B C D) -> (A B C)

CONC: (A B), (C D) ■+ (A B C D)

CONCT: (A B), (C D), (E F) + (A B C D E F)

PREF: (B C), (A B C D) -> (A) LPrefix of (B C) in (A B C D)]

SUFF: (B C), (A B C D) ■+ (D) l Suffix of (B C) in (A B C D) j

and a control structure using COND and NULL.

2*~ Notion_of_tYp_e

A type is a set of lists which can be defined by rules which are summarised as follows [6 1: a) the set of known inputs "x" and the set of outputs "f(x)^M" of the function f to be synthesized are types.

b) if X is a type and f a LISP function, then the set of outputs of f restricted to X as input is a type.

c) if Y is a type and g a LISP function then the set X of x such as g(x) C- Y is a type.

3«~ Segment^].i_on_pattern

Let f be a function to be synthesized and (x, f(x)) an example of "input-output" of this function.

SISP/1 uses a general heuristic to create an expression of the function:

a) segmentation of strings x and y = f(x) into three consecutive segments such that:

CONCT (px, c, sx) ■+ X

CONCT (py, c, sy) -> y

where c denotes the larger string common to x and y, px and py denote the prefixes of c in x and y, sx and sy denote the suffixes of c in x and y.

b) building of relations between these segments.

A "Segmentation Pattern" of (x,y), for all x and y, is defined as the network shown in figure 1.

We can see on this network:

- seven nodes representing types respectively associated to the strings x,y,c,px,py,sx,sy.

- twelve relations between nodes. Each relation consists of a function and a scheme (I,» I«, ...>

I -> J) which indicates the input nodes I , I , ..., I_n in this order and the output node J. This order is represented on the network by a double arrow.

Note that functions FX, FY, GPX, GPY, GSX, GSY are built by SISP/1 using the basic functions LCAR, CDR, LRAC, RDC, and the composition rule. They are chosen of the less possible complexity (the smallest number of basic functions).

In some cases, the segmentation pattern is simpler:

- when one or several strings are empty (NIL),

.. the associated nodes are suppressed from the pattern.

- when two strings are equal, the associated

nodes are joined together. For instance, if x and y are the same, the pattern is reduced to one node; if x and y have no common part, the pattern is reduced to only two nodes.

4*"" \$Y.<lt>£sis fron}_on£_exam£le

The synthesis consists of three steps:

a) SISP/1 generates a network (called a "Segmentation Structure") by the following process:

(1) Generate the segmentation pattern of (x,y).The generation gives the two sets of pairs: {(px» py), (c,py), (sx, py)} Kpx,sy), (c,sy), (sx, sy)}

(2) As long as py and sy are not empty, choose one pair in each set by a heuristic way; for each of these pairs, rename it as (x,y) and go to step 1.

b) SISP/1 looks at the segmentation structure for a lattice in which the minimal and final nodes are respectively X and Y (that is x and y types). This lattice is stepwise constructed using Algorithm 1, defined as follows:

Definiti ons:

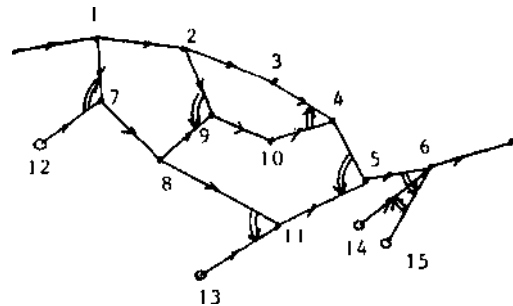
- LAT is the constructed part of the lattice at any step (except in the final step, LAT is not a lattice).

- an incomplete node of LAT is a node such that the relation ending at this node (in LAT) owns some entries which are not connected to X. These nodes are called unsatisfied entries.

- BEG (Z) is the set of nodes in LAT which are less than Z and which are not unsatisfied entries.

- P is a "path" from BEG (U) to V, where U and V are nodes of LAT, if P is an oriented path starting from one node belonging to BEG (U) and ending at V. This path may contain incomplete nodes together with their unsatisfied entries*

Example of LAT:



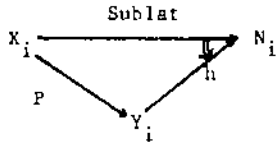
Nodes 6, 7, 11 are incomplete nodes
Nodes 12, 13, 14, 15 are unsatisfied entries
All others nodes are complete nodes.
BEG (9) - {X, 1, 2, 7, 8}

Algorithm 1:

1. LAT «- X
2. Look for a path P between X and Y.
3. Add path P to LAT.
4. If there is no incomplete node in LAT then stop
else select the minimal one and call it N.
(It can be demonstrated that Algorithm 1 generates a set of incomplete nodes which is totally ordered on LAT).

5. Let Y_i be one unsatisfied entry of node N_i .
6. Look i for a path P between $BEG(N_i)$ and Y_i . Let X_i be the origin of P on $BEG(N_i)$ and try to detect a recursivity between X_i and Y_i using Algorithm 2.
7. Go to step 3.

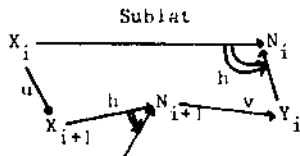
It follows from algorithm 1 that when Algorithm 2 is called, a part of LAT has the following structure:



where sublat is restricted to be a lattice, and Y_i is the previous unsatisfied node.

Algorithm 2:

1. If no path from X_i to N_i (in sublat) matches a subpath of P (in the sense of an identical sequence of operators) then stop algorithm 2, else let $X_{i+1} \rightarrow N_{i+1}$ be the subpath of P which has been matched. The above structure is changed to:



2. If the segmentation structure does not contain a lattice starting in X_{i+1} , ending in N_{i+1} and analogous to sublat then stop algorithm 2 else assume the following recursion:

$$X_i \xrightarrow{\varphi} N_i$$

$$\varphi(x) = h [\text{sublat}(x), v(\varphi(u(x)))]$$

or $X_i \xrightarrow{\varphi} N_i$

$\varphi(x) = h [v(\varphi(u(x))), \text{sublat}(x)]$ depending on the order of arguments of h .

3. Find a primitive stop condition of the recursive function φ as follows: match the operators of path P from X_{i+1} to N_{i+1} in the segmentation structure, then from X_{i+2} to N_{i+2} and so on, until it fails. Assume it fails from X_j to N_j . Find a path w from X_j to N_j . The primitive stop condition is assumed to be:

$$\text{if } X \subset X_j, \text{ then } w(X)$$

4. Reduce the primitive stop condition as follows:

- a) remove Sublat from LAT.
- b) if a relation from N_k to N_{k+1} for every k , $i \leq k < j$, cannot be found in the segmentation structure, then set $P = (X_i \xrightarrow{\varphi} N_i)$ and stop algorithm 2 else let r be the found relation.
- c) find k , the smaller non negative integer such that $w_k = r^k(w(x))$ is not a fixed point of equation:

$$h [\text{sublat}(u^k(x)), v(r(w_k(x)))] = w_k(x)$$

for every $x \in X_j$.

d) Set $P = (X_i \xrightarrow{\varphi} N_i)$

$$\text{with } \varphi(x) = \begin{cases} \text{If } x \in u^k(X_j) \text{ then } z_k(x) \\ \text{else } h [\text{sublat}(x), v(\varphi(u(x)))] \end{cases}$$

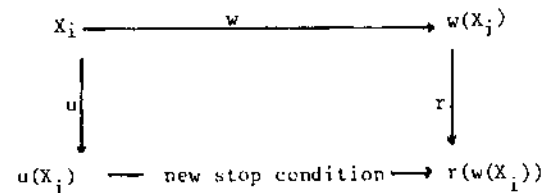
where z_k is solution of the following equation:

$$z_k(u^k(x)) = r^k(w(x)) \text{ for every } x \in X_j$$

Remark : To reduce the primitive stop condition, the following process is iterated:

Suppose the last stop condition is if $x \in X_j$ then $w(x)$

Using the functions u and r , it is possible to calculate $\varphi(x)$, $x \in X_j$, as shown on the following figure:



$$\text{that is: } \varphi(x) = h [\text{sublat}(x), v(\varphi(u(x)))]$$

$$= h [\text{sublat}(x), v(r(w(x)))] = w(x)$$

which means, if the correct answer $w(x)$ is obtained that $w(x)$ is a fixed point of the last equation.

It thus follows that there exists a function z such that:

$$z(u(x)) = r(w(x)) \text{ for every } x \in X_j$$

Thus now the stop condition is:

$$\text{if } x \in u(X_j) \text{ then } z(x)$$

5.- Synthesis from two examples

Let us suppose that after a first example, SISP/1 generates a function which fails on a second example. Let the two examples be (x,y) and (x',y') .

The principle is always to build a structure from the generation of segmentation-patterns; SISP/1

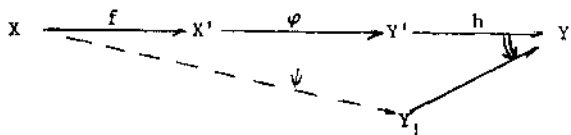
here generates the segmentation patterns associated with the initial pairs $(x,x'),(y,y'),(x,y),(x',y')$ and goes on in the same way as in the first method.

SISP/1 then tries to find a three parts splitted path from X to Y :

- a path from X to X' .
- the function φ itself from X' to Y' .
- a path from Y' to Y .

Remarks: - using this technique, SISP/1 looks explicitly for a recursive form of the function φ .

- when unsatisfied nodes are remaining in the path, SISP/1 generates sub-problems which are to be solved either by algorithm 1 or again by using one more example when algorithm 1 fails [3 bis]. For instance, let Y_i be a remaining unsatisfied node in the path:



SISP provides the following expression of φ :

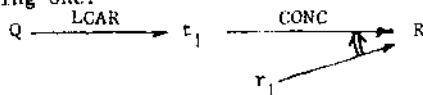
$$\varphi(x) = \begin{cases} \text{if } x \in X' \text{ then } w(x) \\ \text{else } h[\varphi(f(x)), \psi(x)] \end{cases}$$

where ψ is a sub-problem to be solved by SISP/1 and where $w(x)$ is the function which has been found by algorithm 1 working on only the example (x', y') . This stop condition can then be reduced as explained in algorithm 2.

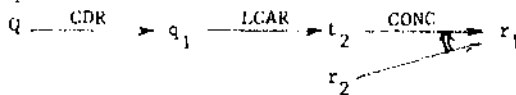
III. PRACTICAL EXAMPLES

1. Let us use our method to find the REVERSE function. The input (A B C D E) is given to SISP/1: it does not know the answer and asks the Professor who returns: (E D C B A). SISP analyses input and output and generates the segmentation structure indicated in figure 3.

SISP looks then for a path from the Question Q containing the list (A B C D E) to the answer R containing (E D C B A) and finds the following one:



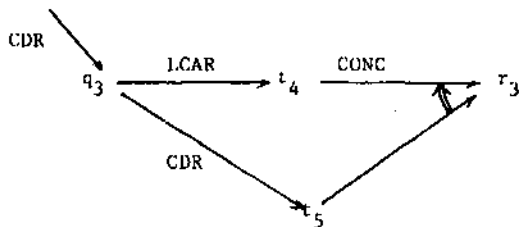
Looking for the unsatisfied entries, SISP finds r_1 . It looks again for a path from $BEG(R)$ to r_1 , and finds the following one:



SISP now examines both paths. The mapping (LCAR-CONC) of the first one matches into the mapping (CDR-LCAR-CONC) of the second one. This statement is sufficient to infer a recursive expression φ :

$$\varphi(x) = \text{CONC}[\varphi(\text{CDR}(x)), \text{LCAR}(x)]$$

SISP has to still find the stop condition. Matching the three operators CDR, LCAR, CONC with the structure, it remarks that it can apply CDR on type q_3 giving t_5 but cannot apply LCAR on type t_5 . SISP/1 thus knows that the stop condition is to be found in this part of the structure:



SISP/1 tries now to find a new binding of the path from t_5 to the unsatisfied entry of r_3 which happens here to be identical to t_5 . It finds the trivial one and generates the primitive stop condition:

$$\text{if } x \in t_5 \text{ then } x$$

SISP now has to reduce the stop condition, using the following mappings:

$$q_i \xrightarrow{\text{CDR}} q_{i+1} \quad r_i \xrightarrow{\text{RDC}} r_{i+1}$$

Assuming that $\text{CDR}((E)) = \text{NIL}$ $\xrightarrow{\varphi}$ $\text{RDC}((E)) = \text{NIL}$ is the first reduced stop condition, SISP calculates now $\varphi(x)$ for every x belonging to t_5 :

$$\begin{aligned} \varphi(E) &= \text{CONC}[\varphi(\text{CDR}((E))), \text{LCAR}((E))] \\ &= \text{CONC}[\varphi(\text{NIL}), (E)] \end{aligned}$$

using the new stop condition: $\varphi(E) = \text{CONC}[\text{NIL}, (E)] = (E)$ which here gives the correct answer.

The process cannot be performed further, because $\text{CDR}(\text{CDR}((E)))$ does not exist. The function generated by SISP is thus:

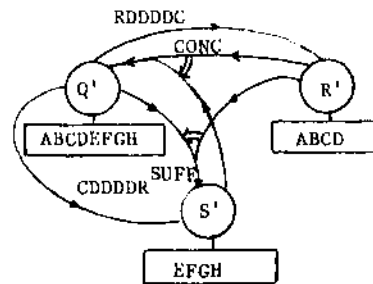
$$\varphi(x) = \begin{cases} \text{if } x = \text{NIL} \text{ then } \text{NIL} \\ \text{else } \text{CONC}(\varphi(\text{CDR}(x)), \text{LCAR}(x)) \end{cases}$$

that is the usual REVERSE.

2. The second example we display now needs more material than the first one since two couples (input - output) are necessary. Let HALF be the function to synthesize:

- first couple: (A B C D E F G H) \rightarrow (A B C D)

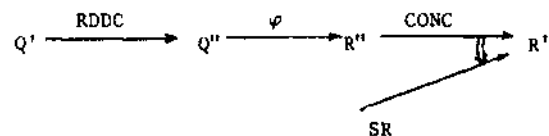
From this first example, SISP/1 constructs the following structure:



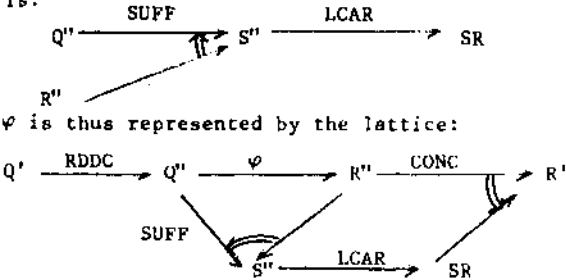
The relation $\varphi: Q' \rightarrow R'$ found here is thus $\varphi(x) = \text{RDDDDC}(x)$.

The professor gives now the following input: (A B C D E F). SISP uses φ to answer (A B), which is false. The professor then gives the correct answer: (A B C D E F) \rightarrow (A B C).

SISP then generates the structure displayed on figure 3. Assuming that R' can be obtained from Q'' using the correct function φ to be synthesized, SISP, as explained before, uses the following path from Q' to R' :



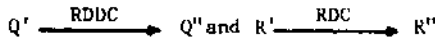
The problem is now to find a path from BEG (R') to SR. The simplest one which is found here is:



φ is thus represented by the lattice:
 $\varphi(x) = \text{CONC} [\varphi(\text{RDDC}(x)), \text{LCAR}(\text{SUFF} [\varphi(\text{RDDC}(x)), \text{RDDC}(x)])]$ with the trivial stop condition:

if $x \in Q'$ then $\text{RDDDC}(x)$

SISP now has to reduce this stop condition using the following mappings:



Assuming that $\text{RDDC}((A B C D E F)) = (A B C D)$

$\varphi \text{ RDC}((A B C D)) = (A B)$
 SISP computes now $\varphi(x)$ for every x belonging to Q'' :
 $\varphi((A B C D E F)) = \text{CONC} [\varphi(\text{RDDC}((A B C D E F))), \text{LCAR}(\text{SUFF} [\varphi(\text{RDDC}((A B C D E F))), \text{RDDC}((A B C D E F))])] = \text{CONC} [\varphi((A B C D)), \text{LCAR}(\text{SUFF} [\varphi((A B C D)), (A B C D)])] = \text{CONC} [(A B), \text{LCAR}(\text{SUFF} [(A B), (A B C D)])] = \text{CONC} [(A B), (C)] = (A B C)$
 which is the correct answer.

The primitive stop condition can thus be reduced to:

if $x \in \text{RDDC}(Q')$ then $\text{RDDC}(x)$

where RDDC is the solution of the following equation on z

$$z(\text{RDDC}(x)) = \text{RDC}(\text{RDDDC}(x)) \text{ for every } x \in Q''.$$

This process is recursively applied and stops when $\text{RDDC}((A B)) = \text{NIL}$. At this step, we obtain the function HALF defined as follows:

$$\varphi(x) = \begin{cases} \text{if } \text{RDDC}(x) = \text{NIL} \text{ then } \text{RDC}(x) \\ \text{else } \text{CONC} [\varphi(\text{RDDC}(x)), \text{LCAR}(\text{SUFF} [\varphi(\text{RDDC}(x)), \text{RDDC}(x)])] \end{cases}$$

IV. LIMITS AND PROSPECTS OF THE METHOD

1.- Prospects

SISP/1 is already able to synthesize most of the functions given in SUMMERS [5] and HEDRICKS [2] in particular it synthesizes the following ones by using algorithm 1:

- (A B C D E) → (E D C B A)
 - (A B C D E) → (A X B X C X D X E X)
 - (A B C D E) → (A A B B C C D D E E)
 - (A B C D E) → (A A B A B C A B C D A B C D E)
- By using two or more examples it synthesizes:
- (A B C D E F G) → (A A G G B B F F C C E E D D)
 - (A B C D E F G H) → (A B C D)
 - (A B C D E) → (E D C B A E D C B E D C E D E)
 - (A B C D) → (D C B A C B A B A A D C B C B B D C C D)

- (A B C D E) → (A B B C C C D D D D E E E E E E)
- (A B C D E F G H) → (D C B A H G F E)
- (A B) → (A A A A A A A A) (cube of the entry length)
- (A B C D) → (A A A A A A A A) (half square): such a way is not always easy to use, as we shall see now:

- assume SISP has to synthesize the HALF function using the previous example (A B C D E F G H) → (A B C D). Algorithm 1 fails and the professor gives as second example:

$$(B C D E F G) \rightarrow (B C D)$$

SISP here generates the following HALF function:

$$\varphi(x) = \begin{cases} \text{if } x = \text{NIL} \text{ then } \text{NIL} \\ \text{else } \text{CONC} [\text{LCAR}(x), \varphi(\text{CDR}(x))] \end{cases}$$

which is much simpler than previous HALF function. This simplicity was however found by the professor who gave better examples.

- assume now that SISP has to synthesize a function using the example

(A B C D E) → (A B B C C C D D D D E E E E E E). Algorithm 1 fails and the professor gives as second example : (A B C D) → (A B B C C C D D D D).

SISP generates the following functions:

$$\varphi(x) = \begin{cases} \text{if } x \in Q'' \text{ then } R'' \\ \text{else } \text{CONC} [\varphi(\text{RDC}(x)), \psi(x)] \end{cases}$$

where $\psi(x)$ is bound to the following subproblem:

$$(A B C D E) \rightarrow (E E E E E)$$

Algorithm 1 fails then to provide a correct function ψ and the professor now has to give the two particularly well chosen examples:

- (B C D E) → (B C C D D D E E E E)
- (B C D) → (B C C D D D)

they allow SISP to generate a new appropriate example in order to synthesize a correct ψ :

$$\psi(x) = \begin{cases} \text{if } x = \text{NIL} \text{ then } \text{NIL} \\ \text{else } \text{CONC} [\psi(\text{CDR}(x)), \text{LRAC}(x)] \end{cases}$$

the stop condition of φ is then found:

if $x = \text{NIL}$ then NIL

the generated function Φ will thus be given by the linear recursive system Φ

$$\Phi: \begin{cases} \varphi(x) = \text{if } x = \text{NIL} \text{ then } \text{NIL} \text{ else } \text{CONC} [\varphi(\text{RDC}(x)), \psi(x)] \\ \psi(x) = \text{if } x = \text{NIL} \text{ then } \text{NIL} \text{ else } \text{CONC} [\psi(\text{CDR}(x)), \text{LRAC}(x)] \end{cases}$$

These two last examples show the main importance of good examples. We hope however that it would be possible to use "bad examples" joined together with a unification process, in order to improve the given "bad examples".

2.- Limits

- with the exception of stop-condition, the functions generated by SISP do not use predicates in their definition. Thus the function:

if "length of x is even" then reverse (x)
else x

cannot be synthesized by SISP. This problem is attacked in [6].

- SISP/1 requires a good sequence of example in order to use the second technique. They have to be of decreasing length and consecutive.
- SISP/1 only works on atomic lists.

V. CONCLUSION

In summary, the described method consists of constructing a structure from an adequate set of

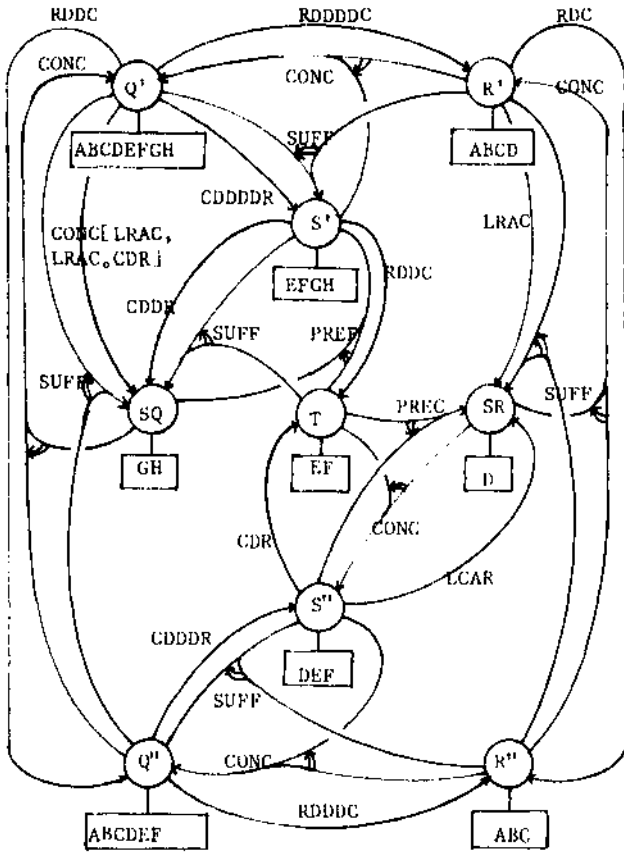


Figure 3: Structure associated to the HALF function.