1064

**Bell Laboratories**　　Cover Sheet for Technical Memorandum

Title- **Agen – An Associative Memory Generator**

Date- **September 2, 1975**

TM- **75-1274-18**

Other Keywords- **Computer Languages**
**Data Management**

| Author | Location | Extension | Charging Case- **39199** |
|---|---|---|---|
| **M. E. Lesk** | **MH 2C-569** | **6377** | Filing Case- **39199-11** |

## ABSTRACT

Agen writes programs to store and retrieve information. Each stored item contains data and a key; to retrieve or allocate the item, just mention the key. Possible applications include symbol tables, sparse arrays, or table lookups of words. Two possible algorithms are built in: the user can request a binary tree search or a hash search form of storage organization.

Agen is a program generator, which translates specifications into a high-level programming language which can be mixed with the user's general purpose code. Each user request for a lookup procedure is translated into a subroutine to perform the search and allocation if needed. This subroutine may be viewed as if it addressed an array of pointers to the user's data items. Agen is available on GCOS and UNIX. It writes lookup subroutines in the C language.

| Pages Text 9 | Other 0 | Total 9 | |
|---|---|---|---|
| No. Figures 1 | No. Tables 1 | No. Refs. 1 | |

DISTRIBUTION
(REFER GEI 13.9-3)

| COMPLETE MEMORANDUM TO | COVER SHEET ONLY TO | COVER SHEET ONLY TO | COVER SHEET ONLY TO | COVER SHEET ONLY TO |
|---|---|---|---|---|
| CORRESPONDENCE FILES | ANDERSON,MS K J | COBEN,ROBERT M | FREEDMAN,M I | HYMAN,B |
|  | ARMSTRONG,DOUGLAS B | COHEN,HARVEY | FREEMAN,K GLENN | IFFLAND,FREDERICK C |
| OFFICIAL FILE COPY | ARNOLD,GEORGE W | COKE,MISS E U | FREEMAN,R DON | IMAGNA,CLYDE P |
| PLUS ONE COPY FOR | ARNOLD,S L | <COLE,LOUIS M | FREIDENREICH,MRS B | IPPOLITI,O D |
| EACH ADDITIONAL FILING | ATAL,B S | COLE,M O | FROST,H BONNELL | IRVINE,M M |
| CASE REFERENCED | BACCASH,MISS J M | COOK,THOMAS J | FULTON,ALAN W | <IVIE,EVAN L |
|  | BADURA,DENNIS C | COOPER,A E | GAJEWSKA,MS HANNA O | JACKOWSKI,D J |
| DATE FILE COPY | BAKER,BRENDA S | COPP,DAVID H | GANNON,T F | JACOBS,H S |
| (FORM E-1328) | BASEIL,RICHARD J | COREY,D A | GARCIA,R F | JAKUBEK,RAYMOND J |
|  | BAUER,MS H A | COSTELLO,PETER E | <GATES,G W | JAMES,DENNIS B |
| 10 REFERENCE COPIES | BAUGH,C R | COSTON,WALTER P | GAWRON,L J | JAMES,J W III |
|  | BAYER,DOUGLAS L | <COULTER,J REGINALD | GAY,FRANCIS A | <JENSEN,PAUL D |
| <AHO,A V | BERNSTEIN,LAWRENCE | CRAGUN,D W | GEARY,M J | JESSOP,WARREN H |
| <BECKER,R A | BERTH,R P | CRANE,RODERICK P | GEPNER,JAMES R | JESSUP,RICHARD F |
| <BOYCE,W M | BEYER,JEAN-DAVID | CRUME,LARRY L | GERGOWITZ,E B | JOHNSTON,WALTER E JR |
| <BROWN,W STANLEY | BICKFORD,N B | CUNNINGHAM,STEPHEN J | GEYLING,F T | JONES,B R |
| <CHAMBERS,J M | BILLINGTON,MISS M J | CUTLER,C CHAPIN | GIBB,KENNETH R | <JUDICE,CHARLES N |
| <COOPER,DENNIS W | BILOWOS,RICHARD M | D ANDREA,MRS LOUISE A | GIBSON,H T JR | KACHURAK,JOSEPH J |
| <FLEISCHER,HERBERT I | BIRCHALL,R H | D STEFAN,D J | <GILLETTE,DEAN | <KAISER,J F |
| <FRASER,A G | BIREN,MRS IRMA B | DAVIS,D R | GIMPEL,JAMES F | KAMINSKY,MICHAEL L |
| <GOLDSTEIN,A JAY | BISHOP,MISS V L | DAVIS,R D | GITHENS,JOHN A | KAPLAN,A E |
| <GUTHERY,SCOTT B | BLINN,JAMES C | DE JAGER,D S | GLASER,W A | KAPLAN,M M |
| <HAMILTON,PATRICIA | BLUE,J L | DESENDORF,JUDITH | GLASSER,ALAN L | KAUFMAN,LARRY S |
| HAMMING,R W | BLUM,MRS MARION | DEUTSCH,DAVID N | GLUCK,F | KAYEL,R G |
| +HANNAY,N B | BLY,JOSEPH A | DICKMAN,B N | <GNANADESIKAN,R | KEANE,E R |
| <JOHNSON,STEPHEN C | BOCKUS,R J | DIMMICK,JAMES O | GOGUEN,MS NANCY | KELLY,L J |
| <KEESE,W M | BODEN,F J | DOLAN,MRS MARIE T | GOLABEK,MISS R | KENNEDY,ROBERT A |
| <KERNIGHAN,BRIAN W | BONANNI,LORENZO E | DOLOTTA,T A | GOLDSMITH,L D | KERTZ,DENIS R |
| <LUDERER,GOTTFRIED W R | BORKIN,S A | DONOFRIO,L J | GORMAN,JAMES E | KEVORKIAN,DOUGLAS E |
| <MARANZANO,JOSEPH F | BOURNE,STEPHEN R | DOWD,PATRICK G | <GRAHAM,R L | KILLMER,JOHN C JR |
| <MC GILL,ROBERT | BOWERS,J L | DRAKE,MRS L | GRAVEMAN,R F | KLEIN,MISS R L |
| MC ILROY,M DOUGLAS | BOWYER,L RAY | DRUMMOND,R E | GREENBAUM,H J | KNOWLTON,KENNETH |
| +MCDONALD,H S | BRADLEY,R H | DUDLEY,MRS E H | GREENSPAN,S J | KNUDSEN,DONALD B |
| MORGAN,S P | BRANDT,RICHARD B | DUFFY,FRANCIS P | GROSS,ARTHUR G | KNUDSON,DEAN O |
| <OSSANNA,J F JR | BROWN,COLIN W | EDELSON,D | GRUCHOWSKI,M P | KORNEGAY,R L |
| PINSON,ELLIOT N | BROWN,WILLIAM R | EDMUNDS,T W | GUERRIERO,JOSEPH R | KREIDER,DANIEL M |
| +PRIM,ROBERT C | BUCHSBAUM,S J | EIGEN,D | HAFER,E H | LA CAVA,JOSEPH L |
| <RALEIGH,THOMAS M | BULLEY,RAYMOND M | EITELBACH,DAVID L | HAIGHT,R C | LAMBERT,MRS C A |
| <ROBERTS,CHARLES S | BURG,F M | ELLIOTT,R J | HALE,A L | LAYTON,RICHARD L |
| TERRY,M E | BURNETTE,W A | ELY,T C | HALL,ANDREW D JR | LEE,DENNIS F |
| <THOMSON,M-L | BURROWS,T A | ERDLE,K W | HALL,MILTON S JR | LEGENHAUSEN,S |
| 29 NAMES | BYRNE,EDWARD R | ERRICHIELLO,PHILIP M | HALL,W G | LESK,MICHAEL E |
|  | BZOWY,D E | ESSERMAN,ALAN R | HALPIN,THOMAS | LESSEK,PETER V |
|  | CAMPBELL,J H | ESTOCK,R G | HAMILTON,MRS L | LEVINE,P R |
| COVER SHEET ONLY TO | CAMPBELL,STEPHEN T | FABISCH,MICHAEL P | HARKNESS,C J | LICWINKO,J S |
|  | CANADAY,RUDD H | FEDER,J | HARRISON,NEAL T | LIEBERT,THOMAS A |
|  | CARAWAY,R E | FELDMAN,STUART I | HARRIS,CHARLES S | LINDERMAN,J |
| CORRESPONDENCE FILES | CARDOZA,WAYNE M | FELS,ALLEN M | HARUTA,K | LIST,WILLIAM H S |
|  | CARRAN,J H | FERIDUN,K K | HAUSE,A D | LIU,T K |
| 4 COPIES PLUS ONE | CARR,DAVID C | FINUCANE,JOHN J | <HAWKINS,DONALD T | <LOGAN,MRS V G |
| COPY FOR EACH FILING | CASPERS,MRS BARBARA E | FISCHER,H B | HEATH,SIDNEY F III | LOMUTO,N |
| CASE | CAVINESS,JOHN D | FLYNN,MS M L | HELD,RICHARD W | LONG,DAVID W |
|  | CHAFFEE,N F | FONG,KENNETH T | HERGENHAN,C B | LORENC,ANTHONY |
| ABRAHAM,STUART A | CHAMBERS,MRS B C | FORTNEY,V J | HEROLD,JOHN W | LUTZ,KENNETH J |
| ACKERMAN,A F | CHEN,EDWARD | FOUNTOUKIDIS,A | HESS,MILTON S | <LYCKLAMA,HEINZ |
| AHRENS,RAINER B | CHEN,STEPHEN | FOWLKES,E B | HOEHN,MISS MARIE J | LYONS,T G |
| ALBERTS,BARBARA A | CHERRY,MS I L | FOX,PHYLLIS | HOLTMAN,JAMES P | MACHOL,R E JR |
| ALCALAY,DAVID | CHIANG,T C | FOY,J C | HONIG,W L | MACRI,PHILIP P |
| ALLISON,CHARLES E | CHODROW,MARK M | FRANKLIN,DANIEL L | HOYT,WILLIAM F | MADDEN,MRS D M |
| ALT,MISS DOROTHY L | CHRIST,C W JR | FRANK,H G | HUDSON,E T | MAHLER,G R |
| AMOSS,JOHN J | CLAYTON,D P | FRANK,MISS A J | HUNNICUTT,CHARLES F | MALCOLM,J A |
| AMRON,I | CLIFFORD,ROBERT M | FRANK,RUDOLPH J | HUPKA,MRS FLORENCE | MALLOWS,COLIN L |
|  |  |  |  | ... |
|  |  |  |  | 454 TOTAL |

+ NAMED BY AUTHOR    > CITED AS REFERENCE    < REQUESTED BY READER    (NAMES WITHOUT PREFIX
WERE SELECTED USING THE AUTHOR'S SUBJECT OR ORGANIZATIONAL SPECIFICATION AS GIVEN BELOW)

MERCURY SPECIFICATION.................................................................................

COMPLETE MEMO TO:
    127-SUP

COVER SHEET TO:
    12-DIR        13-DIR    127

    COPLGP = COMPUTING/PROGRAMMING LANGUAGES/GENERAL PURPOSE
    UNPL# = UNIX/PROGRAMMING LANGUAGES

--------------------------------------------------------------------------------------------------------

TO GET A COMPLETE COPY:

PLEASE SEND A COMPLETE COPY TO THE ADDRESS SHOWN ON THE
OTHER SIDE
NO ENVELOPE WILL BE NEEDED IF YOU SIMPLY STAPLE THIS COVER
SHEET TO THE COMPLETE COPY.

1. BE SURE YOUR CORRECT ADDRESS IS GIVEN ON THE OTHER SIDE.
2. FOLD THIS SHEET IN HALF WITH THIS SIDE OUT AND STAPLE.
3. CIRCLE THE ADDRESS AT RIGHT.  USE NO ENVELOPE.

IF COPIES ARE NO LONGER AVAILABLE PLEASE FORWARD THIS
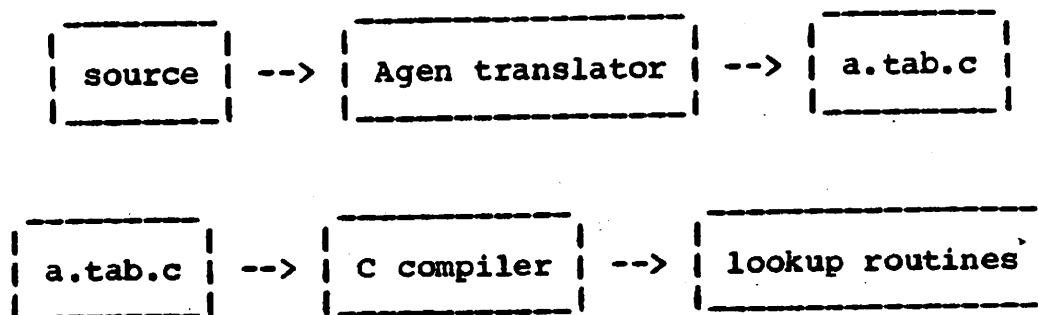REQUEST TO THE CORRESPONDENCE FILES.

MEMORANDUM FOR FILE

Agen is a program generator for table searches. It accepts descriptions of the table keys or indices and writes a subroutine to perform the appropriate lookup. Thus the programmer wishing to store information relating to words in a table can use a character string as a subscript, as if it were an integer.

Agen translates a file of descriptions of lookup routines, called the source file henceforth, into a file of programs implementing the lookups. The generated routines are written on a file named a.tab.c which can be compiled directly or included in a user's C program. The operation of Agen is shown in Figure 1.

Figure 1

Operation of Agen

```
 _____       _____       _____
|           |     |                     |     |           |
|  source   | --> |  Agen translator    | --> |  a.tab.c  |
|_____|     |_____|     |_____|


 _____       _____       _____
|           |     |                     |     |                     |
|  a.tab.c  | --> |    C compiler       | --> |  lookup routines    |
|_____|     |_____|     |_____|
```

Each routine is described by its name, the arguments it expects, and the kind of item stored. The lookup routines return pointers to the stored items. Whenever a new key is used, the storage for the new element appears automatically. Thus a typical source line for Agen might appear:

        word (string) st1;

indicating that the programmer wishes a table lookup function named "word", which has a string as its argument, and returns a pointer to items like "st1". The item stored is

used to declare the size of the table entries, and may be replaced by the number of bytes to be stored. In the user's program, if "st1" is a structure declared like

```
struct {int first, last, count;} st1;
```

then the program may contain operations like

```
word("when")->first = 3;
a = word("in")->last;
if (word("the")->count > 0)
```

just as if there had been a declaration

```
struct {int first, last, count;} word [ ];
```

and "when" or "in" could be magically converted to an appropriate integer subscript.

The kinds of subscripts accepted by Agen programs are integers, characters, and strings. Thus a function can be declared

```
matrix (int, int, int) item;
```

to indicate three integer subscripts. Why is this better than

```
matrix[ a ][ b ][ c ]
```

or some such normal multidimensional array? The difference is that the user need not know the values of a,b, and c; and the amount of storage used is proportional to the amount actually stored into, not the amount declared. In exchange, the fetch and store operations are much slower with Agen than with the normal hardware subscripting operations.

The operation of the Agen lookup routines is easy to describe. First, the arguments are converted to a single key. The key is the linear concatenation of all the <u>int</u> or <u>char</u> subscripts (all are received by the called program as <u>int</u> because of C conventions) followed by the concatenation of all <u>string</u> arguments (separated by '|' characters). This single key, which uniquely identifies the subscripts given, is now used in either a binary tree search or a hash search to find the appropriate element. If no such element exists, it is created, the storage assigned to it cleared, and a pointer to the new element returned. In the normal case, the lookup routine always returns a pointer to a valid item, whether previously defined or not. There can never be more than one item for a given key, since an element is created only if it can not be found in the table.

Sometimes this is not adequate. Thus, in addition to

subscript type arguments, an inquiry type argument may also be used. This is specified as <u>exist</u> in an Agen source line. Such an argument is zero for normal operation; if it is supplied as 1 in a call, it indicates that the caller is merely inquiring as to the existence of an item under the requested subscripts. If the item exists, the usual pointer is returned. If it does not, zero is returned. Thus a definition

        word (string, exist) item;

if used with a call of the form

        word ("course", 0)

creates a new item is nothing is stored under the key "course" but the call

        word ("of", 1);

returns zero if "of" is not defined. This sort of argument has two other uses. If it is -1, it indicates deletion; that is,

        word("human", -1);

would delete an item stored as "human", if any existed. If the <u>exist</u> argument is 2, it indicates "find the next element". In this case, an element different from the one specified is returned, in such a way that continued requests for the "next" element will give the contents of the table. The pointer returned in this case is the pointer to the key, not the pointer to the value. An ordinary search will retrieve the value if desired. If the argument is a simple string or integer, and if the search is a binary tree search, note that the definition of the key is such as to return the table in lexicographic order. When an <u>exist</u> argument of 2 is given and the end of the table is reached, zero is returned.

Table 1 summarizes the uses of the "exist" argument.

There is one other kind of argument that can be supplied to an Agen lookup routine, called "size". Normally, all items stored by a lookup routine are the same size, and the size is given in the declaration. Often, however, it is desirable to store different sized items and specify every time a new item is created how big it is. In this case, a <u>size</u> argument must be declared to the lookup routine and given at the first use of each subscript. Changing the declaration of the "word" routine to

## Table 1

_____

### The EXIST Argument

| Value | Meaning |
|-------|---------|
| -1 | Delete this entry if it exists. Return 0. |
| 0 | Normal operation; create if non-existent. Return a pointer to the data item. |
| 1 | Test existence. Do not create. Return a pointer to the data item if it exists, otherwise return 0. |
| 2 | Find another item. Return a pointer to the next data item (never the one specified by this key). If "no more" items, return 0. |

_____

        word (string, exist, size);

permits calls of the form

        word ("events", 0, 20);

indicating that if an item must be created for "events", it
is to be 20 bytes long. Any existing item will be returned,
regardless of length. This requires the maximum data size
needed for a given key to be specified the first time the
key is used, unless the item is deleted and recreated.

        The exist argument may be set to 2 to find the contents
of a table. Successive calls will return different keys,
which may then be used to obtain the data. With a binary
tree search, "next" really means "next greater than ..." and
thus the table can be obtaned in order of the keys. Note,
however, the description above (page 2) of the way in which
the keys are constructed from the arguments. A single in-
tegral or string key works as expected, but more complicated
keys have the order of their constituents permuted to group
integral and string subscripts. With hash searching, "next"
means "next" in hash table order. Anything missing from the
hash table returns the first item in the table as "next".
Thus, no matter what the form of a table, the entries can be
enumerated by successive probes with an exist argument of 2
as follows:

1) use the null string or zero as a key;
2) use the returned key as the next key;
3) stop when zero is returned.

There are some interactions between the exist and size arguments if both are specified. If the exist argument is 1 the size argument is never relevant, since nothing can be created. If the exist argument is -1, the size argument is essential and must give the size of the element to be deleted.

The source lines defining lookup routines also select the choice of a binary or a hash search. This is specified between the function name and its arguments. The user may give the name of the search method, and/or an estimated numerical length of the table. If the method specified is "hash", the length must be given. The default methods are binary tree search if the length is unknown and hash search if an estimated length was given. In general, hash searches are faster and use less storage. To define "word" as a hash-searched function, write

        word 1500 (string, exist, size);

or

        word hash 1500 (string, exist , size);

The default is the same as

        word binary (string, exist, size);

The general syntax of Agen input is:

        {source lines}
        %%
        {extra material}

Anything after a line containing just %% is copied to the output unchanged. So is anything contained between lines of %{ and %} as in Yacc and Lex. The format of a source line is

        {routine name} {lookup type} {arguments} {item}

where lookup type may be missing, or may be the word "binary" or the word "hash", and may contain a length as an integer. If the word "hash" is included, the length must be given. The arguments are enclosed in parentheses or brackets and separated by commas. Each argument is one of the words int, char, string, exist, or size. The first three indicate subscripts to be used in specifying keys. The size and exist arguments control table handling and allocation.

Arguments may be given in any order. Only one instance of exist or size may appear; any number of subscripts can be used, in any order. The item indicates the size of the items (it can be omitted if there is a size argument) and should either be a number (taken to be in bytes) or the name of a variable to be inserted as an argument to the sizeof() function. Note that the variable has to be declared somewhere; this can either be done directly in the Agen source by use of %{ and %} to bracket the declaration, or in the file which contains a

```
# include "a.tab.c"
```

statement.

As an example, consider the problem of writing a program which reads a text and prints out a table showing the words in the text, the number of occurrences of each word, and the sequence number of the first occurrence. This is very similar to what would be needed to create a symbol table for a computer program, except that the data stored for each symbol would be different. Two source files will be required, one for Agen and one for Lex. Lex [1] is used to generate the program to isolate the words in the input stream. Also, the Lex routine name "yywrap" is used to print the table at the end of the input.

Agen source:

```
%{
        struct {int first, count;} datac;
%}
word (string, exist) datac;
%%
```

**Lex source:**

```
            int kw 0;
%%
[a-zA-Z]+  {
            ++kw;
            word(yytext, 0)->count++;
            if (word(yytext, 0)->count == 1)
                    word(yytext, 0)->first = kw;

            }
            :
\n          ;
%%
yywrap()
{
char *s;
printf("Total First Word\n");
for (s = word("", 2); s != 0; s= word(s,2))
        printf("%5d %5d %s\n",word(s,0)->count,
                word(s,0)->first, s);
printf("total number of words was %d\n",kw);
return(1);
}
# include "a.tab.c"
```

The Agen source defines a table subscripted by character strings (words in this case) and noting for each word the point of first occurrence and the total number of oc-currences. The Lex source uses a rule [a-zA-Z]+ to isolate strings of letters; each such word is looked up and its count incremented. If the count is one, the current word number kw (incremented at each word) is stored in the _first_ field. When the end of input is reached, Lex calls yywrap; this routine runs through the table and prints the statis-tics for each word. Note the procedure for scanning the table: the string s is initialized to the element found by taking the "next" item after the null string. Then succes-sive "next" items are found until the end of the table is reached. Since the search method is a binary tree search, and the keys are simply single character strings, the "next" items are in alphabetical order. If a hash search had been used, the items would be in random order and would have to be sorted on output.

The appropriate commands to compile and run this are

```
agen agensource
lex lexsource
cc lex.yy.c -ll -lp
a.out <input
```

and if the input file is

when in the course of human events, it becomes
necessary for one people to dissolve the
political bands which have connected
them with another, and to assume among the
powers of the earth the separate and equal
station to which

The output is:

| Total | First | Word |
|-------|-------|------|
| 1 | 28 | among |
| 2 | 25 | and |
| 1 | 24 | another |
| 1 | 27 | assume |
| 1 | 18 | bands |
| 1 | 9 | becomes |
| 1 | 21 | connected |
| 1 | 4 | course |
| 1 | 15 | dissolve |
| 1 | 33 | earth |
| 1 | 37 | equal |
| 1 | 7 | events |
| 1 | 11 | for |
| 1 | 20 | have |
| 1 | 6 | human |
| 1 | 2 | in |
| 1 | 8 | it |
| 1 | 10 | necessary |
| 2 | 5 | of |
| 1 | 12 | one |
| 1 | 13 | people |
| 1 | 17 | political |
| 1 | 30 | powers |
| 1 | 35 | separate |
| 1 | 38 | station |
| 5 | 3 | the |
| 1 | 22 | them |
| 3 | 14 | to |
| 1 | 1 | when |
| 2 | 19 | which |
| 1 | 23 | with |

total number of words was 40

As an example of timings, the above program was run
under several options. As above, with a binary tree search
and the portable library used for input and output, it
processes the Declaration of Independence at 116 words per
processor second. With the redundant searches eliminated by
remembering a pointer to the just-made search, the speed is
160 words per processor second. By using a hash search and
UNIX-tailored input and output routines, the speed can be
raised to about 350 words per processor second.

In general, the hash searches are faster and require less memory than the binary tree searches. The main reason for using a binary tree search is that the user may be unwilling or unable to estimate the table size in advance. Hash tables are not grown; the size must not exceed 50% more than the original estimate. A single hash search takes somewhat over half a millisecond; a single binary tree search perhaps 30-50% more time. Even more unfortunate is that the storage for each item includes: a) the memory consumed by the item itself; b) the memory consumed by the key for the item; c) four extra words of memory for the binary tree search or two extra words for the hash search; and d) breakage in the allocator of two words in the case of hash searches and three words in the case of binary searches. When the items and keys are short, there is a lot of wasted storage.

Agen is available on GCOS as well as on UNIX. The speeds of GCOS programs are comparable to those of UNIX programs. To access Agen on GCOS, say

./agen sourcefile

and the a.tab.c file is written as on UNIX.

There exist plans to improve Agen, by producing lookup routines more carefully tailored to the user's applications. This should improve the speed of search routines produced by Agen.

M. E. Lesk

MH-1274-MEL-unix

## References

[1] Lex - A Lexical Analyzer Generator, M. E. Lesk, 1975 TM 75-1274-15.