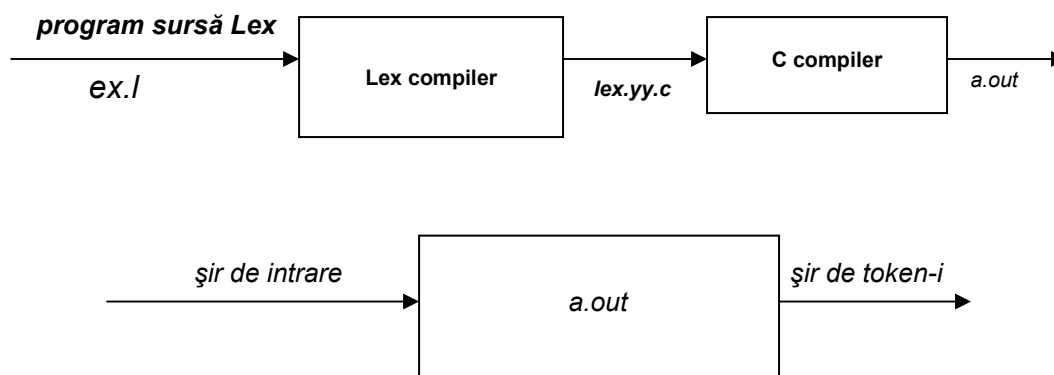


## Limbajul LEX (FLEX)

Limbajul LEX a fost introdus de Lesk și Schmidt în 1975, ca un limbaj pentru specificarea analizorului lexical.

În general Lex-ul este utilizat pentru a genera analizorul lexical, **a.out**, în modul următor:



## Specificații Lex (Flex)

Un program Lex are trei părți:

```
declarații
% %
reguli de traducere
% %
proceduri auxiliare
```

Secțiunea de declarații include declarații de variabile, constante și definiții regulate. Declarațiile de variabile și constante și eventual alte instrucțiuni se scriu în limbaj C și trebuie incluse între semnele `%{` și `%}`.

Definițiile regulate sunt similare cu cele prezentate înainte și sunt folosite drept componente ale expresiilor regulate care apar în regulile de traducere.

Regulile de traducere ale unui program Lex sunt instrucțiuni de forma:

```
p1 {acțiune 1}
p2 {acțiune 2}
.....
pn {acțiune n}
```

unde fiecare  $p_i$  este o expresie regulată și fiecare acțiune  $i$  este un fragment de program care descrie ce acțiune trebuie executată când forma  $p_i$  produce o lexemă. În Lex, acțiunile sunt scrise în limbajul C.

În secțiunea a III-a apar procedurile auxiliare care sunt necesare acțiunilor. Aceste proceduri pot fi compilate separat și adăugate apoi la analizorul lexical.

### Exemplu 1:

Un program LEX simplu, care transformă literele șirurilor de caracter introduse în litere mari.

```
%%  
[a-z]          {putchar(toupper(yytext[0]));}  
%%  
main()  
{  
    yylex();  
}
```

### Compilări și execuție:

Dacă numele programului de mai sus este `ex1.1` atunci trebuie să executăm următoarele comenzi:

```
flex ex1.1
```

Se generează fișierul `lex.yy.c`.

```
gcc -o ex1 lex.yy.c -lfl
```

Compilăm fișierul C generat, `lex.yy.c`, generăm executabilul cu numele `ex1` cu opțiunea `-o` iar cu opțiunea `-lfl` realizăm link cu anumite librării.

Executăm fișierul executabil generat:

```
./ex1
```

sau

```
./ex1<in.txt
```

### Reguli Lex

O regulă Lex este formată dintr-o expresie regulată și o acțiune asociată, formată din instrucțiuni C. În momentul în care o expresia regulată este recunoscută pe parcursul citirii caracterelor de către scanner, se execută acțiunile corespunzătoare se execută.

În specificarea expresiilor regulate Lex suportă o serie de convenții, notații (clase de caractere, repetări specifice, etc.).

### Exemplu

```
%{
    static unsigned nrNum = 0;
}%
cifra      0|1|2|3|4|5|6|7|8|9
numar      {cifra}+
%%
{numar}    nrNum++; /*număr câte numere naturale apar în text*/
```

Dacă pentru un *pattern* am mai multe instrucțiuni ce reprezintă acțiunea, acestea se plasează într-un bloc de instrucțiuni, între { }.

### Observații:

- Tot ce este scris în prima secțiune a declarațiilor între semnele %{ și %} și tot ce se află în ultima secțiune este copiat identic în programul C rezultat din sursa Lex.
- Secțiunile 1 și 3 sunt opționale
- Regulile din secțiunea 2 trebuie să respecte:
  - încep din prima coloană (nu se lasă spații albe sau tab-uri) în care se scrie *pattern*-ul, dat de o expresie regulată
  - expresiile regulate conțin caractere text (din alfabetul de intrare) și caractere operator
  - *pattern*-urile sunt separate de acțiuni prin spații albe sau tab-uri.

## Expresii regulate în LEX

Expresiile regulate sunt formate din *caractere text* și din *caractere operator*.

Caracterele text: literele alfabetului și cifrele.

Caracterele operator: „ \ [ ] ^ - ? . \* + | ( ) \$ / { } % < >

### Observație:

Caracterele operator pot fi transformate în caractere text, dacă sunt precedate de caracterul \ sau prin punerea între ghilimele.

## Notații în LEX

### Clase de caractere:

Sunt specificate cu ajutorul parantezelor pătrate [ ].

**Exemple:** [abc] reprezintă a sau b sau c  
[0-9] înseamnă 0|1|2|3|4|5|6|7|8|9

**[0-9A-Za-z]** înseamnă orice literă sau cifră

În afară de caracterele \, - și ^, majoritatea celorlalte caractere operator își pierd semnificația și rămân simple caractere text.

**Caracterul – (minus):** se folosește pentru stabilirea unui domeniu:

[0-9] = cifrele de la 0 la 9

Dacă dorim utilizarea caracterului – în interiorul domeniului cu semnificația de caracter text atunci trebuie introdus chiar la început sau precedat de caracterul \.

Caracterul \ are aceeași semnificație de operator descris mai sus.

### Caracterele operator:

Caracterul ^                      utilizat la începutul unei clase de caractere indică complementarea mulțimii specificate. De exemplu, **[^0-9]** înseamnă orice caracter diferit de cifră.

**^x** recunoaște **x** doar la început de linie.

Caracterul .                      înseamnă orice caracter cu excepția caracterului \n (newline)

Caracterul \$                      **x\$** recunoaște **x** doar la sfârșit de linie.

Caracterul ?                      reprezintă opționalitate: **ab?c** înseamnă șirul abc sau ac.

Caracterul \*                      reprezintă stelarea (închiderea Kleene) ca la expresii regulate.

Caracterul +                      una sau mai multe instanțe (ca la expresii regulate). **a+** este echivalent cu **aa\***.

Caracterul |                      reprezintă „sau” logic.

Parantezele ( )                      se folosesc pentru grupare: exemplu: (ab|cd)?(ef)\*

Parantezele { }                      pot fi folosite pentru desemnarea unui nume din definițiile regulate **{numar}**

sau astfel:

**a{3}** reprezintă șirul **aaa** (de 3 ori a)

**x{3,7}** reprezintă **x** între **3** și **7** ori.

**[a-z]{1,5}** reprezintă orice șir format din litere mici de lungime **1** până la **5**

Operatorul / se folosește pentru indicarea contextului drept: astfel **ab/cd** semnifică faptul că **ab** este acceptat doar dacă urmează **cd**

## Variabile globale Lex

Funcția de citire a tokenilor **yylex** nu are nici un argument și returnează o valoare întreagă. Alte variabile globale care oferă mai multe informații despre tokenii citați sunt:

- **yytext** este un șir de caractere null-terminat care conține caracterele lexemei recunoscute. Această variabilă este declarată în fișierul **lex.yy.c**.
- **yylen** este un întreg care memorează lungimea lexemei memorate în **yytext**. Această variabilă este declarată de asemenea în fișierul **lex.yy.c**.
- **yyval** este o variabilă globală folosită pentru a stoca atributele tokenilor. Această variabilă este de tipul **YYSTYPE** și este de obicei un union de diferite câmpuri pentru tokeni diferiți.
- **yyloc** este o variabilă globală care stochează locația, linia și coloana token-ului. Este de tipul **YYLTYPE**.

## Acțiuni Lex

- Există o acțiune predefinită: dacă programul detectează un șir de caractere care nu se potrivește cu nici un *pattern*, atunci se execută acțiunea predefinită care constă din copierea la ieșirea standard a șirului respectiv. De aceea, dacă se dorește evitarea acestui lucru, trebuie definite acțiuni pentru orice șir de caractere posibil.
- Regula care nu efectuează nimic: dacă vreau ca un șir de caractere să fie pur și simplu ignorat nu pun ca regulă nimic, doar ;

**Exemplu:** ignorarea spațiilor albe

```
[ \t\n] ;
```

- Rutina **yymore( )**: are ca rezultat concatenarea următoarei lexeme de cea păstrată în **yytext**.

**Exemplu:** dacă definesc tokenul șir de caractere ca fiind orice caractere puse între „”, în care pot să apară și „”, dar doar precedate de \.

```
\"[^"]*" {
    if(yytext[yylen-1]=='\\') yymore();
    else printf(„sirul de caractere =%s\n”,yytext);
}
```

## Rezolvarea ambiguităților

Limbajul Lex poate trata seturi de reguli ambigue. Anume atunci când mai multe reguli sunt aplicabile Lex alege una dintre ele astfel:

1. Cel mai lung șir care se potrivește pe un *pattern* este ales.

2. Dacă pentru cel mai lung șir care se potrivește pe un *pattern* există mai mult de o regulă atunci se alege prima întâlnită.

### Exemplu:

```
[0-9]+          printf(„numar intreg”);  
[0-9]+(\\.[0-9]+)?  printf(„numar fractionar”);
```

Dacă întâlnește 923.1 va alege a doua regulă.

### Atenție:

Dacă inversez ordinea regulilor va alege și pentru orice număr întreg varianta numărului fracționar, deoarece se potrivește pe ambele reguli și atunci o alege pe prima.

### Exemplu 2:

```
%{  
    int numchar = 0, numword = 0, numline = 0;  
}%  
%%  
\\n          {numline++; numchar+=2;}  
[ ^ \\t\\n ]+ {numword++; numchar+= yyleng;}  
 .           {numchar++;}  
%%  
main()  
{  
    yylex();  
    printf(„Numarul de caractere = %d\\n”, numchar);  
    printf(„Numarul de cuvinte = %d\\n”, numword);  
    printf(„Numarul de linii = %d\\n”, numline);  
}
```

### Compilări și execuție:

```
flex ex2.1
```

```
gcc -o ex2 lex.yy.c -lfl
```

```
./ex2
```

sau

```
./ex2<in.txt
```

### Tema de laborator 1:

Scrieți un program Lex (Flex) care să recunoască următoarele șiruri de caractere peste alfabetul  $V = \{0, 1\}$ :

- Toate cuvintele care se termină în 01
- Toate cuvintele care conțin exact un 0
- Toate cuvintele care conțin un număr par de 1 și nici un 0
- Toate cuvintele care conțin un număr par de 1
- Toate cuvintele care conțin subșirul 01
- Toate cuvintele care nu conțin subșirul 01
- Toate cuvintele care nu conțin subșirul 011

### Tema de laborator 2:

- Să se transforme toate secvențele de caractere albe într-un singur caracter alb și secvențele de newline într-un singur newline dintr-un fișier text obținând rezultatul în alt fișier text.
- Să se determine toți identificatorii, toate numerele întregi și toate numerele reale dintr-un fișier text.
- Să se determine șirurile de caractere dintr-un fișier text ce conține un program în limbajul C.
- Să se elimine comentariile dintr-un fișier text ce conține un program în limbajul C.

### Tema de laborator 3:

Completați exemplul pentru analizorul lexical din man pages corespunzător comenzii flex.

```
/* scanner for a toy Pascal-like language */
%{
    /* need this for the call to atof() below */
    #include <math.h>
}%
DIGIT      [0-9]
ID          [a-z][a-z0-9]*
%%
{DIGIT}+    {
                printf( "An integer: %s (%d)\n", yytext,
                        atoi( yytext ) );
            }

{DIGIT}+"."{DIGIT}*    {
                        printf( "A float: %s (%g)\n", yytext,
                                atof( yytext ) );
                    }

if|then|begin|end|procedure|function    {
```

```

                                printf( "A keyword: %s\n", yytext );
                                }

{ID}          printf( "An identifier: %s\n", yytext );

"+"|"-"|"*"|"/"      printf( "An operator: %s\n", yytext );

"{ "[^]\n}*"          /* eat up one-line comments */

[ \t\n]+              /* eat up whitespace */

.                    printf( "Unrecognized character: %s\n", yytext );

%%
main( argc, argv )
    int argc;
    char **argv;
{
    ++argv, --argc; /* skip over program name */
    if ( argc > 0 )
        yyin = fopen( argv[0], "r" );
    else
        yyin = stdin;
    yylex();
}

```

## **Bibliografie:**

Marinescu D, „Tehnici de compilare”  
 Marinescu D, „Limbae formale și teoria automatelor”