# Exercise P1. Lexical Scanner for Simplified C Programming Language

## 1 Aim of the Exercise

The aim of the exercise is to build a simple scanner for a much simplified version of the programming language C. The task of the analyzer is:

- to recognize tokens of C language and to determine their values

- to remove blanks and comments

- to recognize given directives

- to recognize lexical errors

## 2 Preliminaries

After turning on the computer, one should select Linux, and log in as *student*. One should open a console window (e.g. press `Alt-F2` and type `xterm`), create one's own directory using a command `mkdir` *family name of the user*, and a subdirectory for the current exercise. Download files for the C language from the Moodle web page of the course for the subject *Lexical Analysis*. The following files are to be found there:

- `p1c.pdf` — manual (just being read)

- `Makefile` — needed for compilation with the command `make`

- `common.h` — header file defining the greatest length of strings

- `c.l` — skeletal lexical analyzer that needs to be completed; take a closer look at the definition of `process_token()`, which should be used in the rules

- `c.y` — parser that is needed only for declaring tokens and for invoking the lexical analyzer

- `test1.c` — correct test program

- `test2.c` — test program with errors that should be detected

After having completed the exercise, the directory should be removed.

## 3 Tasks

The supplied skeletal lexical analyzer should be extended so that it works correctly on supplied test programs. The analyzer should print information on recognized tokens in three columns:

1. matched text

2. recognized token

3. value of the token (only when it makes sense)

Function `process_token` is designed to print that information. The function returns a recognized token, so an action in a rule recognizing a token should contain **return process_token(. . . )** with appropriate parameters.

The supplied code needs to be completed with the following items **using the ordering given below**:

A. printing one's own name (in the `bison` program)

B. detecting keywords in test programs (defined in the source code for `bison`)

C. removing blanks

D. removing one-line comments

E. recognition of multi-character operators (`<=`, `++`,. . . ) that appear in test programs

F. recognition of identifiers

G. recognition of integers and floating point numbers

H. recognition of strings without start conditions

I. recognition of character constants

J. recognition of one-character tokens: operators and punctuation

K. recognition of include directives

L. recognition of strings using start conditions

M. removal of multi-line comments using start conditions

N. detection of comment end sequence without the beginning sequence using start conditions

O. detection of failure to close a comment with indications of the line where the comment begins

# 4 Grading

All items are graded as 1 point. If needed, items from K to O can be completed **at home for half a point each**. The file developed in the lab should be uploaded before the end of the class on Moodle. **The lexical analyzer will be needed for the next exercise**. Make sure that recognized tokens are handed over to the parser using `return`. Tokens being removed, e.g. comments or white spaces, should not call `return`.

# 5 Start Conditions

- Start condition active at the start of the program: `INITIAL`

- Declaraction: `%x condition1, condition2,...`

- Matching in a start condition:
  ```
  <con1> re1      action1;
  <con1,con2,INITIAL>re2 action2;
  <*>re3          action3
  ```

- changing start condition: `BEGIN condition4`

- current start condition: `YYSTATE`

- checking the current start condition after all input data has been read: in function `yywrap`, which must be defined, and which must return 1

# 6 Test Data — File test1.c

```c
/****************************************************/
/* Program ASCII - wyswietla rozszerzone kody ASCII */
/****************************************************/
#include <stdio.h>
#include "test.h"
unsigned char uc; // zmienna sterujaca petli typu char
int fromASCII = 128, toASCII = 255;
long int x[10];
void main( void )
{
  struct data {
    int rok;
    int miesiac;
    int dzien;
  };
  data poczatek, koniec;
  int i;
  printf( "Rozszerzone kody ASCII\n\n");
  for ( uc = fromASCII; uc <= toASCII; uc1++ ) {
    printf( "%3d:%2c", uc, uc ); printf("\n");
```

```
21      }
22      int x1 = fromASCII + 2 * ( 20 +  toASCII );   /* int */
23      double realTest = 12.34e-12 + .56 + 78.; /* double */
24      x[0] = 1;
25      for (i = 1; i < 10; i++) {
26         x[i] = x[i-1] * i * i;
27      }
28      poczatek.rok = 2018;
29      poczatek.miesiac = 10;
30      poczatek.dzien = 1;
31 }
```

## 7  Test Data — File test2.c

```
1  /*****************************************************/
2  /* Program ASCII - wyswietla rozszerzone kody ASCII */
3  /*****************************************************/
4  #include <stdio.h>
5  #include "test.h"
6  unsigned char uc; // zmienna sterujaca petli
7  int fromASCII = 128, toASCII = 255;
8  void main( void )
9  {
10          printf("\n\n\nRozszerzone kody ASCII\n\n");
11          for (uc = fromASCII; uc <= toASCII; uc1++)
12          {
13                  printf("%3d:%2c", uc, uc);
14          }
15 }
16 int x1 = fromASCII + 2 * ( 20 +  toASCII ); /* te linie /* sluza
17 * / wylacznie celom testowym ;-) */
18 double realTest = 12.34 + .56 + 78.;
19 */ // nieotwarty komentarz
20 "Niezamknieta stala tekstowa
21 /* niezamkniety komentarz
```

## 8  Output of the Lexical Analyzer for test1.c

```
1  Author: First name and family name
2  yytext                Symbol type       Symbol value as string
3
4  Przetwarzanie dyrektywy #include <stdio.h>
5  Przetwarzanie dyrektywy #include "test.h"
6  unsigned              KW_UNSIGNED
7  char                  KW_CHAR
8  uc                    IDENT             uc
9  ;                     ;
10 int                   KW_INT
11 fromASCII             IDENT             fromASCII
12 =                     =
13 128                   INTEGER_CONST     128
14 ,                     ,
15 toASCII               IDENT             toASCII
16 =                     =
17 255                   INTEGER_CONST     255
18 ;                     ;
19 long                  KW_LONG
20 int                   KW_INT
21 x                     IDENT             x
22 [                     [
23 10                    INTEGER_CONST     10
```

```
24  ]                          ]
25  ;                          ;
26  void                       KW_VOID
27  main                       IDENT             main
28  (                          (
29  void                       KW_VOID
30  )                          )
31  {                          {
32  struct                     IDENT             struct
33  data                       IDENT             data
34  {                          {
35  int                        KW_INT
36  rok                        IDENT             rok
37  ;                          ;
38  int                        KW_INT
39  miesiac                    IDENT             miesiac
40  ;                          ;
41  int                        KW_INT
42  dzien                      IDENT             dzien
43  ;                          ;
44  }                          }
45  ;                          ;
46  data                       IDENT             data
47  poczatek                   IDENT             poczatek
48  ,                          ,
49  koniec                     IDENT             koniec
50  ;                          ;
51  int                        KW_INT
52  i                          IDENT             i
53  ;                          ;
54  printf                     IDENT             printf
55  (                          (
56  "Rozszerzone kody AS       STRING_CONST      "Rozszerzone kody ASCII\n\n"
57  )                          )
58  ;                          ;
59  for                        KW_FOR
60  (                          (
61  uc                         IDENT             uc
62  =                          =
63  fromASCII                  IDENT             fromASCII
64  ;                          ;
65  uc                         IDENT             uc
66  <=                         LE
67  toASCII                    IDENT             toASCII
68  ;                          ;
69  uc1                        IDENT             uc1
70  ++                         INC
71  )                          )
72  {                          {
73  printf                     IDENT             printf
74  (                          (
75  "%3d:%2c"                  STRING_CONST      "%3d:%2c"
76  ,                          ,
77  uc                         IDENT             uc
78  ,                          ,
79  uc                         IDENT             uc
80  )                          )
81  ;                          ;
82  printf                     IDENT             printf
83  (                          (
84  "\n"                       STRING_CONST      "\n"
85  )                          )
86  ;                          ;
87  }                          }
88  int                        KW_INT
89  x1                         IDENT             x1
```

| | | |
|---|---|---|
| 90 | = | = |
| 91 | fromASCII | IDENT | fromASCII |
| 92 | + | + |
| 93 | 2 | INTEGER_CONST | 2 |
| 94 | * | * |
| 95 | ( | ( |
| 96 | 20 | INTEGER_CONST | 20 |
| 97 | + | + |
| 98 | toASCII | IDENT | toASCII |
| 99 | ) | ) |
| 100 | ; | ; |
| 101 | double | KW_DOUBLE |
| 102 | realTest | IDENT | realTest |
| 103 | = | = |
| 104 | 12.34e−12 | FLOAT_CONST | 12.34e−12 |
| 105 | + | + |
| 106 | .56 | FLOAT_CONST | .56 |
| 107 | + | + |
| 108 | 78. | FLOAT_CONST | 78. |
| 109 | ; | ; |
| 110 | x | IDENT | x |
| 111 | [ | [ |
| 112 | 0 | INTEGER_CONST | 0 |
| 113 | ] | ] |
| 114 | = | = |
| 115 | 1 | INTEGER_CONST | 1 |
| 116 | ; | ; |
| 117 | for | KW_FOR |
| 118 | ( | ( |
| 119 | i | IDENT | i |
| 120 | = | = |
| 121 | 1 | INTEGER_CONST | 1 |
| 122 | ; | ; |
| 123 | i | IDENT | i |
| 124 | < | < |
| 125 | 10 | INTEGER_CONST | 10 |
| 126 | ; | ; |
| 127 | i | IDENT | i |
| 128 | ++ | INC |
| 129 | ) | ) |
| 130 | { | { |
| 131 | x | IDENT | x |
| 132 | [ | [ |
| 133 | i | IDENT | i |
| 134 | ] | ] |
| 135 | = | = |
| 136 | x | IDENT | x |
| 137 | [ | [ |
| 138 | i | IDENT | i |
| 139 | − | − |
| 140 | 1 | INTEGER_CONST | 1 |
| 141 | ] | ] |
| 142 | * | * |
| 143 | i | IDENT | i |
| 144 | * | * |
| 145 | i | IDENT | i |
| 146 | ; | ; |
| 147 | } | } |
| 148 | poczatek | IDENT | poczatek |
| 149 | . | . |
| 150 | rok | IDENT | rok |
| 151 | = | = |
| 152 | 2018 | INTEGER_CONST | 2018 |
| 153 | ; | ; |
| 154 | poczatek | IDENT | poczatek |
| 155 | . | . |

```
156  miesiac              IDENT              miesiac
157  =                    =
158  10                   INTEGER_CONST      10
159  ;                    ;
160  poczatek             IDENT              poczatek
161  .                    .
162  dzien                IDENT              dzien
163  =                    =
164  1                    INTEGER_CONST      1
165  ;                    ;
166  }                    }
```

## 9  Output of the Lexical Analyzer for test2.c

```
 1  Author: First name and family name
 2  yytext               Symbol type        Symbol value as string
 3
 4  Processing directive #include <stdio.h>
 5  Processing directive #include "test.h"
 6  unsigned             KW_UNSIGNED
 7  char                 KW_CHAR
 8  uc                   IDENT              uc
 9  ;                    ;
10  int                  KW_INT
11  fromASCII            IDENT              fromASCII
12  =                    =
13  128                  INTEGER_CONST      128
14  ,                    ,
15  toASCII              IDENT              toASCII
16  =                    =
17  255                  INTEGER_CONST      255
18  ;                    ;
19  void                 KW_VOID
20  main                 IDENT              main
21  (                    (
22  void                 KW_VOID
23  )                    )
24  {                    {
25  printf               IDENT              printf
26  (                    (
27  "\n\n\nRozszerzone k STRING_CONST       "\n\n\nRozszerzone kody ASCII\n\n"
28  )                    )
29  ;                    ;
30  for                  KW_FOR
31  (                    (
32  uc                   IDENT              uc
33  =                    =
34  fromASCII            IDENT              fromASCII
35  ;                    ;
36  uc                   IDENT              uc
37  <=                   LE
38  toASCII              IDENT              toASCII
39  ;                    ;
40  uc1                  IDENT              uc1
41  ++                   INC
42  )                    )
43  {                    {
44  printf               IDENT              printf
45  (                    (
46  "%3d:%2c"            STRING_CONST       "%3d:%2c"
47  ,                    ,
48  uc                   IDENT              uc
49  ,                    ,
50  uc                   IDENT              uc
```

6

| | | | |
|---|---|---|---|
| 51 | ) | ) | |
| 52 | ; | ; | |
| 53 | } | } | |
| 54 | } | } | |
| 55 | int | KW_INT | |
| 56 | x1 | IDENT | x1 |
| 57 | = | = | |
| 58 | fromASCII | IDENT | fromASCII |
| 59 | + | + | |
| 60 | 2 | INTEGER_CONST | 2 |
| 61 | * | * | |
| 62 | ( | ( | |
| 63 | 20 | INTEGER_CONST | 20 |
| 64 | + | + | |
| 65 | toASCII | IDENT | toASCII |
| 66 | ) | ) | |
| 67 | ; | ; | |
| 68 | double | KW_DOUBLE | |
| 69 | realTest | IDENT | realTest |
| 70 | = | = | |
| 71 | 12.34 | FLOAT_CONST | 12.34 |
| 72 | + | + | |
| 73 | .56 | FLOAT_CONST | .56 |
| 74 | + | + | |
| 75 | 78. | FLOAT_CONST | 78. |
| 76 | ; | ; | |
| 77 | Unexpected closure of a comment in line 19 | | |
| 78 | String opened in line 20 not closed in the same line | | |
| 79 | Missing comment closure for comment opened in line 21 | | |