

CS480

Translators

Finishing Lex Analysis

Chap. 3

Lexical-Analyzer Generator

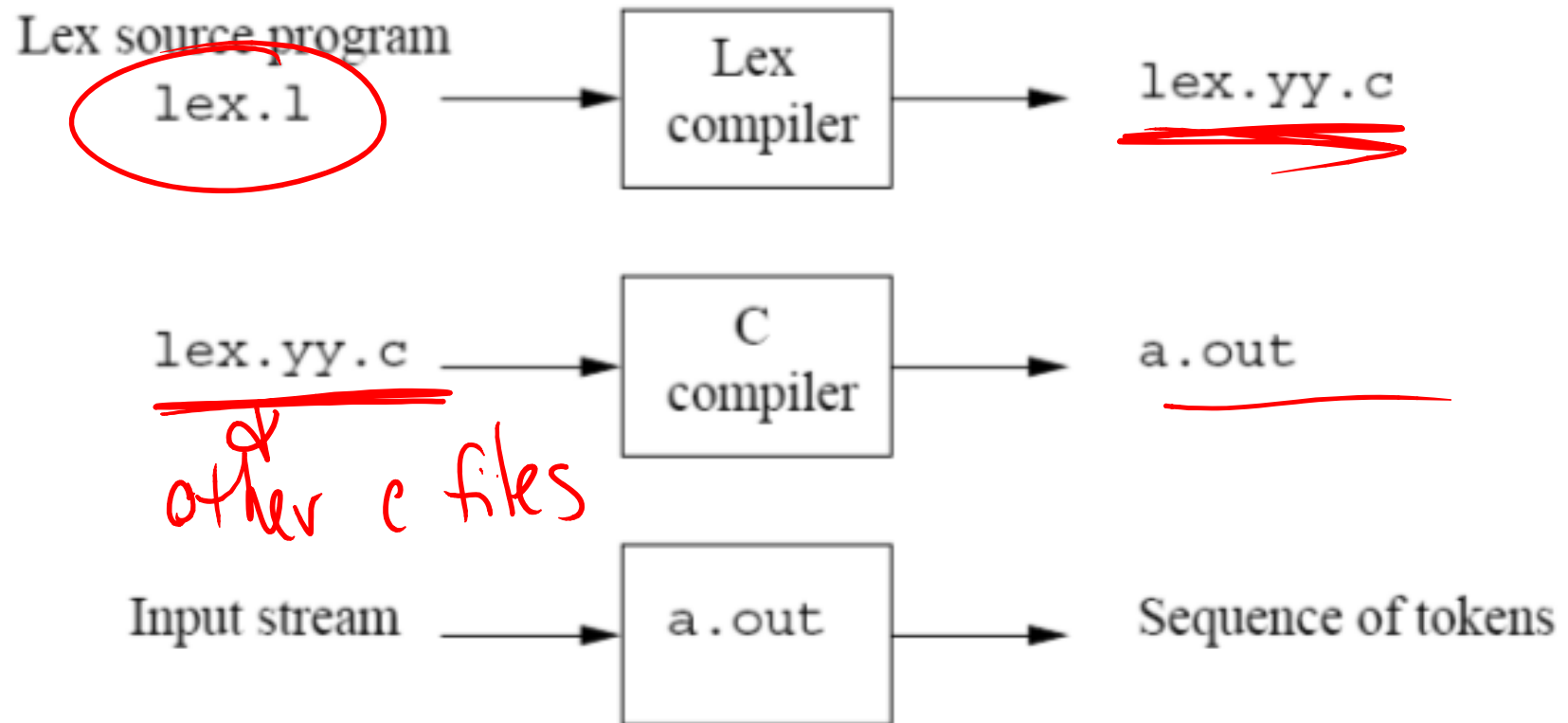


Figure 3.22: Creating a lexical analyzer with Lex

What does Lex do?

- A program that takes regular expressions and action pairs, and builds a recognizer.
- Why actions?

Using Lex

- 1. LEX makes NO assumptions about tokens, must recognize ALL of input, including spaces, comments even illegal symbols ('.' default action).
- 2. Multiple rules can match the same input. Ambiguity is resolved using the following two rules:
 - A) choose pattern which matches the longest input string
 - B) if two patterns match the same size string, choose that which was listed first.

```
%{
/* definitions of manifest constants
LT, LE, EQ, NE, GT, GE,
IF, THEN, ELSE, ID, NUMBER, RELOP */
```

defines for coding

```
%}

/* regular definitions */
delim    [ \t\n]
ws        {delim}+
letter    [A-Za-z]
digit     [0-9]
id         {letter}({letter}|{digit})*
number     {digit}+(\.{digit}+)?(E[+-]?{digit}+)?
```

reg ex
for what
tokens
are built
upon the
token

```
%%

{ws}      { /* no action and no return */ }
if         { return(IF); }
then       { return(THEN); }
else       { return(ELSE); }
{id}       { yylval = (int) installID(); return(ID); }
{number}   { yylval = (int) installNum(); return(NUMBER); }
"<"       { yylval = LT; return(RELOP); }
"<="      { yylval = LE; return(RELOP); }
"="        { yylval = EQ; return(RELOP); }
"<>"      { yylval = NE; return(RELOP); }
">"        { yylval = GT; return(RELOP); }
">="      { yylval = GE; return(RELOP); }
```

tokens

%%

```
int installID() { /* function to install the lexeme, whose
                  first character is pointed to by yytext,
                  and whose length is yyleng, into the
                  symbol table and return a pointer
                  thereto */
}

int installNum() { /* similar to installID, but puts numer-
                   ical constants into a separate table */
}
```

Figure 3.23: Lex program for the tokens of Fig. 3.12

```

1  %{
2      #define DIGIT 2
3      #define NUM 3
4  %{
5  ws      [ \t\n]+
6  digit   [0-9]
7  num     {digit}+
8
9  %%
10 {ws}      {} // ignore whitespace
11 {digit}   { return(DIGIT); } // match on digit
12 {num}     { return(NUM); } // match integer numbers
13 .        { yyerror(); } //error on anything else
14 %%
15
16 int main(void)
17 { int n;
18   while ( n = yylex() ) // call scanner until it returns 0 for EOF
19       // token code, lexeme string, length
20       printf ("Code:%d Text:%s TextLength:%d\n", n, yytext, yyleng);
21
22   return 0;
23 }
24
25 int yyerror(void) // default action in case of error in yylex()
26 { printf(" error\n");
27   exit(0);
28 }
29
30 int yywrap(void)
31 { return 1; } // won't compile on Linux w/o it
"lex.1" 31L, 682C written

```

Handwritten notes:

- Red arrows pointing to lines 2, 3, 5, 6, 7, 10, 11, 12, 13, 16, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31.
- Red text: "yylex code that parser needs" with arrows pointing to lines 18 and 19.
- Red text: "tokens" with a bracket pointing to lines 10, 11, 12, 13.
- Red text: "codes" with an arrow pointing to line 16.
- Red text: "default action in case of error in yylex()" with an arrow pointing to line 25.
- Red text: "won't compile on Linux w/o it" with an arrow pointing to line 31.

```
access.engr.orst.edu - PuTTY
flip1 ~/480test 184% lex lex.l
flip1 ~/480test 185% gcc lex.yy.c -o scanner
flip1 ~/480test 186% cat input
56 7 8
flip1 ~/480test 187% ./scanner < input
Code:3 Text:56 TextLength:2
Code:2 Text:7 TextLength:1
Code:2 Text:8 TextLength:1
flip1 ~/480test 188% vi input
flip1 ~/480test 189% cat input
56 7 8 if
flip1 ~/480test 190% ./scanner < input
Code:3 Text:56 TextLength:2
Code:2 Text:7 TextLength:1
Code:2 Text:8 TextLength:1
error
flip1 ~/480test 191% █
```

pass

everything
you can
think of
that should
be handled

fail
fail
fail?


```
1 %{\n2     #define DIGIT 2\n3     #define NUM 3\n4 %}\n5 ws      [ \\t\\n]+\n6 digit   [0-9]\n7 num     {digit}+\n8\n9 %%\n10 {ws}      { }          // ignore whitespace\n11 {digit}   { return(DIGIT); } // match on digit\n12 {num}     { return(NUM); }  // match integer numbers\n13 .         { yyerror(); }   //error on anything else\n14 %%\n15\n16 /*\n17 int main(void)\n18 {\n19     int n;\n20     while ( n = yylex() ) // call scanner until it returns 0 for EOF\n21         // token code, lexeme string, length\n22         printf ("Code:%d Text:%s TextLength:%d\\n", n, yytext, yylen);  
23     return 0;\n24 }\n25 */\n26\n27 int yyerror(void) // default action in case of error in yylex()\n28 {\n29     printf(" error\\n");\n30     exit(0);\n31 }\n32\n33 int yywrap(void)\n34 {\n35     return 1; } // won't compile on Linux w/o it
```

```
1 #include <stdio.h>
2
3 int yylex(); // scanner prototype
4 extern char* yytext;
5 extern int yyleng;
6
7 int main(void)
8 { int n;
9   while ( n = yylex() ) // call scanner until it returns 0 for EOF
10     // output the token code, lexeme string, length
11     printf("Code:%d Lexeme:%s Length:%d\n", n, yytext, yyleng);
12
13   return 0;
14 }
15
```

access.engr.orst.edu - PuTTY

flip1 ~/480test 207% lex lex.l

flip1 ~/480test 208% gcc test.c lex.yy.c -o scanner

flip1 ~/480test 209% ./scanner < input

Code:3 Lexeme:56 Length:2

Code:2 Lexeme:7 Length:1

Code:2 Lexeme:8 Length:1

error

flip1 ~/480test 210% cat input

56 7 8 if

flip1 ~/480test 211% ■

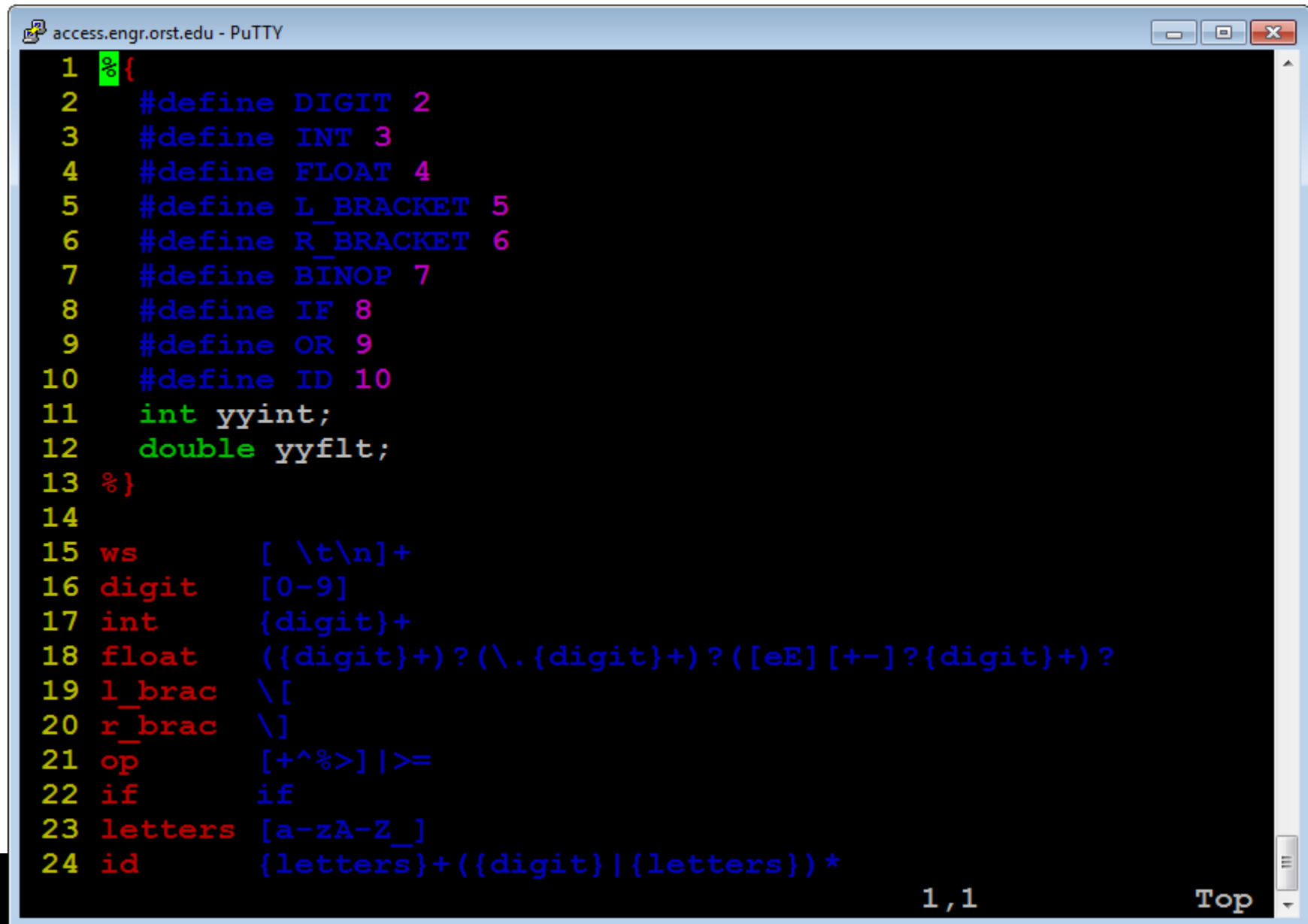
```
1 %{
2     #define DIGIT 2
3     #define NUM 3
4 }%
5 ws      [ \t\n]+
6 digit   [0-9]
7 num     {digit}+
8
9 %%
10 {ws}    { }           // ignore whitespace
11 {num}   { return(NUM); } // match integer numbers
12 {digit} { return(DIGIT); } // match on digit
13        { yyerror(); }  //error on anything else
14 %%
15
16 /*
17 int main(void)
18 {   int n;
19     while ( n = yylex() ) // call scanner until it returns 0 for EOF
20         // token code, lexeme string, length
21         printf ("Code:%d  Text:%s  TextLength:%d\n", n, yytext, yyleng);
22
23     return 0;
24 }
25 */
26
27 int yyerror(void) // default action in case of error in yylex()
28 {   printf(" error\n");
29     exit(0);
30 }
31
32 int yywrap(void)
33 {   return 1; } // won't compile on Linux w/o it
```

```
access.engr.orst.edu - PuTTY
flip2 ~/480test 164% cat input
56 7 8 if
flip2 ~/480test 165% lex lex.l
lex.l:12: warning, rule cannot be matched
flip2 ~/480test 166% gcc test.c lex.yy.c -o scanner
flip2 ~/480test 167% ./scanner < input
Code:3 Lexeme:56 Length:2
Code:3 Lexeme:7 Length:1
Code:3 Lexeme:8 Length:1
error
flip2 ~/480test 168% █
```

1 %{
2 #define DIGIT 2
3 #define NUM 3
4 void *yyval;
5 %}
6 ws [\t\n]+
7 digit [0-9]
8 num {digit}+
9
10 %%
11 {ws} { } // ignore whitespace
12 {num} { yyval=yytext; return(NUM); } // match integer numbers
13 {digit} { return(DIGIT); } // match on digit
14 . { yyerror(); } //error on anything else
15 %%
16
17 /*
18 int main(void)
19 { int n;
20 while (n = yylex()) // call scanner until it returns 0 for EOF
21 // token code, lexeme string, length
22 printf ("Code:%d Text:%s TextLength:%d\n", n, yytext, yyleng);
23
24 return 0;
25 }
26 */
27
28 int yyerror(void) // default action in case of error in yylex()
29 { printf(" error\n");
30 exit(0);
31 }
32
33 int yywrap(void)
34 { return 1; } // won't compile on Linux w/o it


```
access.engr.orst.edu - PuTTY
flip2 ~/480test 174% lex lex.l
lex.l:13: warning, rule cannot be matched
flip2 ~/480test 175% gcc test.c lex.yy.c -o scanner
flip2 ~/480test 176% ./scanner < input
Code:3 Lexeme:56 Length:2
Code:3 Lexeme:7 Length:1
Code:3 Lexeme:8 Length:1
error
flip2 ~/480test 177% █
```


How about for our IBTL...



The image shows a PuTTY terminal window titled "access.engr.orst.edu - PuTTY". The terminal displays a C preprocessor and lexer code. The code is as follows:

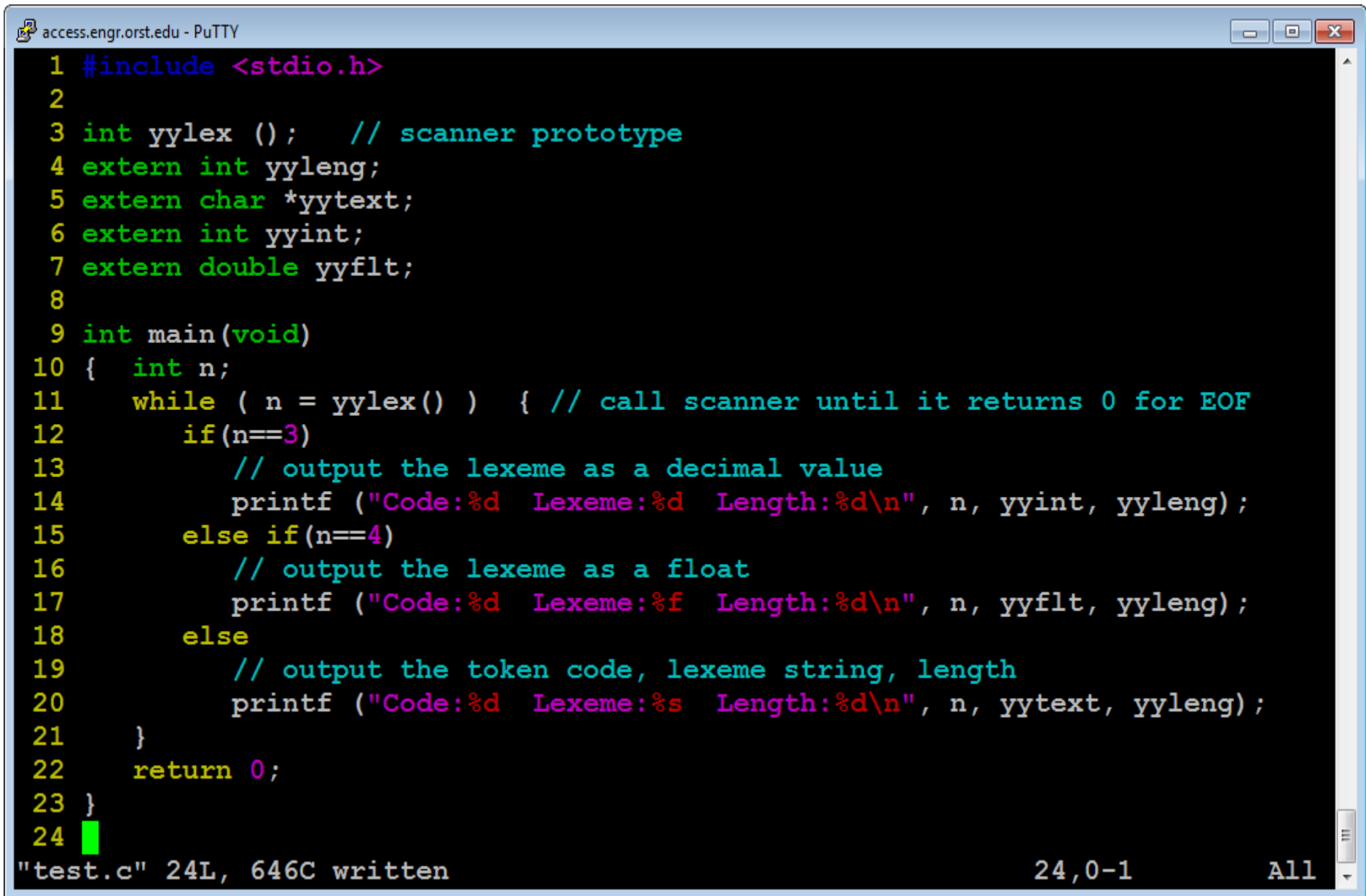
```
1  %{
2      #define DIGIT 2
3      #define INT 3
4      #define FLOAT 4
5      #define L_BRACKET 5
6      #define R_BRACKET 6
7      #define BINOP 7
8      #define IF 8
9      #define OR 9
10     #define ID 10
11     int yyint;
12     double yyflt;
13  %}
14
15  ws      [ \t\n]+
16  digit   [0-9]
17  int     {digit}+
18  float   ({digit}+)?(\.{digit}+)?([eE][+-]?{digit}+)?
19  l_brac  \[
20  r_brac  \]
21  op      [+^%>]|>=
22  if      if
23  letters [a-zA-Z_]
24  id      {letters}+({digit}|{letters})*
```

The terminal window has a status bar at the bottom right showing "1,1" and a "Top" button.

How about for our IBTL...

```
access.engr.orst.edu - PuTTY
21 op      [+^%>]|>=
22 if      if
23 letters [a-zA-Z_]
24 id      {letters}+(({digit})|{letters})*
25
26 %%
27 {ws}     { } // ignore whitespace
28 {int}     { yyint=atoi(yytext); return(INT); } // match integer numbers
29 {float}   { yyflt=atof(yytext); return(FLOAT); } // match on float
30 {l_brac}  { return(L_BRACKET); }
31 {r_brac}  { return(R_BRACKET); }
32 {op}      { return(BINOP); }
33 {if}      { return(IF); } //match if before id or it will be id
34 {id}      { return(ID); }
35 .         { yyerror(); } //error on anything else
36 %%
37
38 int yyerror(void) // default action in case of error in yylex()
39 { printf(" error\n");
40   exit(0);
41 }
42
43 int yywrap(void)
44 { return 1; } // won't compile on Linux w/o it
"lex.l" 44L, 997C written                                     37,0-1 Bot
```

How about for our IBTL...



```
access.engr.orst.edu - PuTTY
1 #include <stdio.h>
2
3 int yylex ();    // scanner prototype
4 extern int yyleng;
5 extern char *yytext;
6 extern int yyint;
7 extern double yyflt;
8
9 int main(void)
10 {   int n;
11     while ( n = yylex() ) { // call scanner until it returns 0 for EOF
12         if(n==3)
13             // output the lexeme as a decimal value
14             printf ("Code:%d Lexeme:%d Length:%d\n", n, yyint, yyleng);
15         else if(n==4)
16             // output the lexeme as a float
17             printf ("Code:%d Lexeme:%f Length:%d\n", n, yyflt, yyleng);
18         else
19             // output the token code, lexeme string, length
20             printf ("Code:%d Lexeme:%s Length:%d\n", n, yytext, yyleng);
21     }
22     return 0;
23 }
24 █
"test.c" 24L, 646C written                24,0-1                All
```