

Appendix B

Cheat-sheet: LEX

B.1 LEX Source

The general format of a Lex source is:

Listing B.1.1: format.l

C code

```
{definitions}
%%
{rules}
%%
{user subroutines}
```

The second `%%` is optional, but the first is required to mark the beginning of the rules. The absolute minimum Lex program is described below and translates into a program which copies the input to the output unchanged.

Listing B.1.2: smallest_lex_file.l

C code

```
%%
```

The **definitions** section may contain:

- translations for repetitive regular expressions
- `includes` of libraries and global variable definitions, enclosed in `%{ %}`
- a list of start conditions
- others, such as:
 - selection of a host language
 - a character set table
 - adjustments to the default size of arrays within Lex itself

The **rules** section represents the user's control decisions. Rules are a table in which:

- the left column contains regular expressions
- the right column contains actions, program fragments to be executed when the expressions are recognized

The **user subroutines** section is often omitted, but may contain overrides for the `yywrap()` function which is called whenever Lex reaches an end-of-file, and any other user defined functions.



Note B.1.1

In the definitions and rules sections, any indented text, or text enclosed in `{% and %}` is copied verbatim to the output (with the `{%}`'s removed). The `{%}`'s must appear unindented on lines, and by themselves.

A sample lex file containing typically used elements:

Listing B.1.3: sample.l

C code

```
%{
#include <stdio.h>
#include "y.tab.h"
int c;
%}
D      [0-9]
%%
" "    ;
[a-z]  {
        c = yytext[0];
        yylval.a = c - 'a';
        return(LETTER);
      }
{D}    {
        c = yytext[0];
        yylval.a = c - '0';
        return(DIGIT);
      }
[~a-z0-9\b] {
        c = yytext[0];
        return(c);
      }
%%
yywrap ()
{
    < ... some post processing ... >
    return 1;
}
```

Concept B.1.1: Comments

In Lex, C-like comments such as `/* ... */` are valid, but for the rules section, only comments for lines that are indented (with whitespaces or tabs) will be correctly identified.

```

%{ /*
This is a valid lex file. The comments are all correct for Lex. This is Comment 0.
Also note the three other comments in various locations in this code.
As always, note the use of white-space before any C code.

Comment 0: Within %{ ...%}.
Comment 1: In the comments section, preceded by a blank.
Comment 2: In the body, within a compound C statement (in braces).
Comment 3: In the body, preceded by a blank. Also note semi-colon.
*/
%}
/*comment1*/
int i,x,y;
D      [0-9]
H      [a-fA-F]
%%
[xX]({H}|{D})+ { /*comment2*/
    y=0;
    for (i=1 ; i < strlen (yytext) ; i++) {
        x=yytext [i];
        if (x >= '0' && x <= '9') {x = x - '0';}
        if (x >= 'A' && x <= 'F') {x = x - 'A' + 10;}
        if (x >= 'a' && x <= 'f') {x = x - 'a' + 10;}
        y = y * 16 + x;}
    printf ("%s\t%d\t%d\n", yytext, strlen (yytext)-1 ,y);}
.* /*comment3*/;

```

B.2 Regular Expressions in LEX

Regular expressions in Lex may specify a set of characters/strings to be matched. It can contain:

- expressions specific to exact letters/digits to be matched
- expressions specific to exact operator characters `" \ [] ^ - ? . * + | ($ / { } .`
- expressions with meta characters that have a special meaning
- expressions with repetitions, choices and other properties

For some of the use cases in the tables below, multiple examples are provided, separated by a semicolon.

Table B.1: General Meta Characters

Usage	Utility	Example	Example Notes
text characters	Any Character corresponding to those in the input string (letters and digits)	<code>a</code> ; <code>123</code>	matches the character <code>a</code> , matches the digits <code>1</code> , <code>2</code> and <code>3</code> in this order
<code>.</code>	Any Character Except New Line	<code>a.*b</code>	matches any string that begins with <code>a</code> , ends with <code>b</code> , and has any other character(s) except for newline inbetween
<code>-</code>	Specifies ranges when used in a character class. If needed as non-special character, it should be specified at the beginning or end of the character class.	<code>[a-z0-9<>]</code> ; <code>[-+0-9]</code>	matches lowercase letters, digits, angular brackets; matches plus and minus signs and digits
<code>{x}</code>	Translation of x from the definitions section	<code>{D}</code> , where D is defined as <code>D [0-9]</code>	uses the translation for the Digits definition <code>D</code> to match digits
<code>\</code>	Escapes characters to avoid special behavior, or specifies a special character	<code>[\40-\176]</code> ; <code>\n</code> ; <code>*</code>	matches all ASCII characters from blank to tilde; matches newline; matches the star character
<code>\n</code>	Matches New Line	<code>[a-z][0-9]\n</code>	matches lines that contain a letter and a digit
<code>\t</code>	Matches Tab	<code>[\t]</code> ;	skips tabs and spaces in the input string

Table B.2: Anchor Meta Characters

Usage	Utility	Example	Example Notes
<code>^</code>	If <code>^</code> is the first character in an expression, then the expression will only be matched at the beginning of a line (used outside of a character class).	<code>^abc</code>	matches the string <code>abc</code> if found at the beginning of a line
<code>^</code>	Negates the set that follows if used within a character class.	<code>[^abc]</code>	matches any character except for <code>a</code> , <code>b</code> and <code>c</code>
<code>\$</code>	If <code>\$</code> is the last character in an expression, then the expression will only be matched at the end of a line.	<code>abc\$</code>	matches the string <code>abc</code> if found at the end of a line
<code>x/y</code>	Trailing Context. If <code>x</code> and <code>y</code> are two LEX regular expressions then <code>x/y</code> is another LEX regular expression that matches <code>x</code> if <code>x</code> is followed by <code>y</code> . After use in this context, <code>y</code> is returned to the input before the action is executed (the action only sees the text matched by <code>x</code>)	<code>ab/cd</code> ; <code>ab/\$</code>	matches the string <code>ab</code> if followed by <code>cd</code> ; equivalent to <code>ab/\n</code>
<code><y>x</code>	<code>x</code> when lex is in start condition <code>y</code>	<code><ONE>x</code>	matches the character <code>x</code> if start condition <code>ONE</code> is met

Table B.3: Grouping Meta Characters

Usage	Utility	Example	Example Notes
<code>[]</code>	Matches Characters specified by the set within brackets	<code>[xy]</code> ; <code>[x-z]</code>	matches characters <code>x</code> or <code>y</code> ; matches characters <code>x</code> , <code>y</code> or <code>z</code>
<code>[^]</code>	Matches Characters NOT specified by the set within the brackets	<code>[^0-9]</code> ; <code>[^abc]</code>	anything but digits; anything but <code>a</code> , <code>b</code> or <code>c</code>
<code>x y</code>	Either <code>x</code> Or <code>y</code>	<code>ab cd</code>	matches the string <code>ab</code> or the string <code>cd</code>
<code>x </code>	The action for <code>x</code> is the action for the next rule	<code>[a-z] </code> followed by <code>[A-Z] ;</code>	does nothing for both lowercase and uppercase letters
<code>()</code>	Groups the expressions specified inside. Not necessary for Either Or, but useful for complex expressions	<code>(ab cd+)?(ef)*</code>	matches strings such as <code>abefef</code> , <code>efefef</code> , <code>cdef</code> , <code>cddd</code> , but not strings like <code>abc</code> , <code>abcd</code> , or <code>abcdef</code>

Table B.4: Quantifier Meta Characters

Usage	Utility	Example	Example Notes
<code>x*</code>	0 or More instances of x	<code>[A-Za-z][A-Za-z0-9]*</code>	all alphanumeric strings that start with a letter
<code>x+</code>	1 or More instances of x	<code>[a-z]+</code>	any string containing lowercase letters except for the empty string
<code>x?</code>	0 or One instance of x	<code>ab?c</code>	matches <code>ac</code> or <code>abc</code>
<code>x{m,n}</code>	m through n occurrences of x	<code>[a-z]{1,3}</code>	any string containing 1 to 3 lowercase letters

**Note B.2.1**

When a special character needs to be used as a text character, it must be escaped using `"\"`. Besides `"\"`, you can also enclose a character within quotes to avoid escaping it.

Example: `xyz ++` is the same as `xyz "++"` and as `" xyz ++"`.