

# Recap

In summary:

- The specification of a lexical analyser generator consists of two parts:
  1. Specification of tokens – done through regular expressions.
  2. Specification of actions - done through action routines.
- The lexical analyser generator:
  - Processes the regular expressions and forms a graph called DFA.
  - Copies the action routines without any change.
  - Adds a driver routine whose behaviour we described.

These three things put together constitutes the lexical analyser.

# Issues

- What are regular expressions? How can they be used to describe tokens?
- How can regular expressions be converted to DFA?

# Introduction to Regular Expressions

A regular expressions denote a set of strings, also called *a language*. For example,  $\mathbf{a^*b}$  denotes the language  $\{\mathbf{b, ab, aab, aaab, \dots}\}$ . We denote the language of a regular expression  $r$  as  $L(r)$ .

A single character is a regular expression.

- Examples:  $\mathbf{a, Z, \backslash n, \backslash t}$ .
- Denotes a singleton set containing the character.  $\mathbf{a}$  denotes the set  $\{\mathbf{a}\}$ .

# Introduction to Regular Expressions

$\epsilon$  is a regular expression.

- Denotes  $\{\epsilon\}$ , the set containing the empty string.

# Introduction to Regular Expressions

If  $r$  and  $s$  are regular expressions then  $r|s$  is a regular expression.

- Examples:  $a|b| \dots |z|A|B| \dots |Z$  and  $0|1| \dots |9$ . Let us call these regular expressions **LETTER** and **DIGIT**.
- $L(r|s)$  is the union of strings in  $L(r)$  and  $L(s)$ .

# Introduction to Regular Expressions

If  $r$  and  $s$  are regular expressions then  $rs$  is a regular expression.

- Examples: `begin` – with an assumed associativity.
- $\{\text{LETTER}\}(\{\text{LETTER}\}|\{\text{DIGIT}\})^*$ .
  - Notice that the **braces** required around `LETTER` is a lex requirement and denotes that it is a **synonym for a regular expression and not the literal LETTER.**
- $L(rs)$  is the concatenation of strings  $x$  and  $y$  such that  $x \in L(r)$  and  $y \in L(s)$ .

# Introduction to Regular Expressions

If  $r$  is a regular expressions then  $r^*$  is a regular expression.

- Examples:  $(\{\text{LETTER}\} \mid \{\text{DIGIT}\})^*$
- $L(r^*)$  is the concatenation of zero or more strings from  $L(r)$ .  
Concatenation of zero strings is defined to be the null string.

# Introduction to Regular Expressions

If  $r$  is a regular expressions then  $(r)$  is a regular expression. Parentheses are used for grouping.

- Examples:  $(\{\text{LETTER}\} | \{\text{DIGIT}\})^*$
- The language denoted by  $(r)$  is  $L(r)$ .



# Introduction to Regular Expressions

Shorthand: If  $r$  is a regular expressions then  $r^+$  is a regular expression.

- Examples:  $\{\text{DIGIT}\}^+$
- $L(r^+)$  is the concatenation of one or more strings from  $L(r)$ .
- $r^+ = rr^*$ .

# Introduction to Regular Expressions

Shorthand: If  $r$  is a regular expressions then  $r?$  is a regular expression.

- Examples:  $\{\text{DIGIT}\}?$  denotes zero or one occurrence of a digit.
- $r?$  stands for zero or one occurrence of strings in  $r$ .
- $r? = \epsilon | r$

# Regular expressions provided by Lex

`c` matches the character `c` literally (e.g., `a` in `apple` matches `a`).

`\c` treats `c` as a special character or escape sequence (e.g., `\*` matches `*` literally, while `*` alone has a special meaning in regex).

<u>Expression</u>	<u>Describes</u>	<u>Example</u>
<code>c</code>	any character <code>c</code>	<code>a</code>
<code>\c</code>	character <code>c</code> literally	<code>\*</code>
<code>"s"</code>	string <code>s</code> literally	<code>"**"</code>
<code>.</code>	any character except newline	<code>a.*b</code>
<code>^</code>	beginning of a line	<code>^abc</code>
<code>\$</code>	end of line	<code>abc\$</code>
<code>[s]</code>	any character in <code>s</code>	<code>[abc]</code>
<code>[^s]</code>	any character not in <code>s</code>	<code>[^abc]</code>
<code>r*</code>	zero or more <code>r</code> 's	<code>a*</code>
<code>r+</code>	one or more <code>r</code> 's	<code>a+</code>
<code>r?</code>	zero or one <code>r</code>	<code>a?</code>
<code>r<sub>1</sub>r<sub>2</sub></code>	<code>r<sub>1</sub></code> then <code>r<sub>2</sub></code>	<code>ab</code>
<code>r<sub>1</sub> r<sub>2</sub></code>	<code>r<sub>1</sub></code> or <code>r<sub>2</sub></code>	<code>a b</code>
<code>(r)</code>	<code>r</code>	<code>(a b)</code>
<code>r<sub>1</sub>/r<sub>2</sub></code>	<code>r<sub>1</sub></code> when followed by <code>r<sub>2</sub></code>	<code>abc/123</code>

The expression `r1/r2` means `r1` when followed by `r2`, but it does not include `r2` in the match.

## Example of token specification in Lex

```
[ \t\n]+          { /*no action, no return*/ }
if                { return(IF); }
then              { return(THEN); }
else              { return(ELSE); }
{letter}({letter}|{digit})* {yylval=install_id(); return(ID);}

-?{digit}+(\.{digit}+)?(E[+-]?{digit}+)?
                                {yylval=atof(yytext); return(NUM);}

"<"                {yylval=LT; return(RELOP);}
"<="              {yylval=LE; return(RELOP);}
"+"               {yylval=PLUS; return(ADDOP);}
"*"               {yylval=MULT; return(MULOP);}
```

# LEXICAL ERRORS

Primarily of two kinds:

1. Lexemes whose length exceed the bound specified by the language.
  - In (old time) Fortran, an identifier more than 7 characters long is a lexical error.
  - Most languages have a bound on the precision of numeric constants. A constant whose length exceeds this bound is a lexical error.
2. Illegal characters in the program.
  - The characters ~, & and @ occurring in a Pascal program (but not within a string or a comment) are lexical errors.
3. Unterminated strings or comments.

# Handling Lexical Errors

The action taken on detection of an error are:

1. Issue an appropriate error message.
2.
  - Error of the first type—the entire lexeme is read and then truncated to the specified length. Generates a warning.
  - Error of the second type—
    - Skip illegal character—this is what was discussed earlier.
    - A possibility which is rarely practiced—pass the character to the parser which has better knowledge of the context in which error has occurred. This opens up more possibilities of recovery - replacement instead of deletion.
  - Error of the third type—wait till end of file and issue error message.