

# Lexical Analysis

CMPT 379: Compilers

Instructor: Anoop Sarkar

[anoopsarkar.github.io/compilers-class](https://anoopsarkar.github.io/compilers-class)

# Regular Languages

- The set of regular languages: each element is a regular language

$$- R = \{R_1, R_2, \dots, R_n, \dots\}$$

- Each regular language is an example of a (formal) language, i.e. a set of strings

$$R_1 = \{a\}, R_2 = \{a, aa, aaa, \dots\}, R_3 = \{b\},$$

$$R_4 = \{ba, ab\}, R_5 = \{\epsilon, b, bb, bbb, \dots\}, \dots$$

# Regular Expressions: Definition

- Meaning function  $L$  maps syntax to semantics
  - $L(r)$  = Meaning of regexp  $r$  (a regular language)
  - $a^* = \{\epsilon, a, aa, aaa, \dots\}$
  - $\epsilon = \text{''}$
  - $c = c$
  - $A|B = A \cup B$
  - $AB = \{ab \mid a \in A \text{ and } b \in B\}$
  - $A^2 = \{xy \mid x \in A \text{ and } y \in A\}$
  - $A^* = A^0 \cup A^1 \cup A^2 \cup A^3 \dots$

# Regular Expressions: Definition

- Meaning function  $L$  maps syntax to semantics
  - $L(r)$  = Meaning of regexp  $r$  (a regular language)
  - $L(a^*) = \{\epsilon, a, aa, aaa, \dots\}$
  - $L(\epsilon) = \{\epsilon\}$
  - $L(c) = \{c\}$
  - $L(A|B) = L(A) \cup L(B)$
  - $L(AB) = \{ab \mid a \in L(A) \text{ and } b \in L(B)\}$
  - $L(A^2) = \{xy \mid x \in L(A) \text{ and } y \in L(A)\}$
  - $L(A^*) = L(A^0) \cup L(A^1) \cup L(A^2) \cup L(A^3) \dots$

# Regular Expressions

- Why use meaning function?
  - Make clear what is syntax and what is semantics
  - Allow us to consider notation as a separate issue

# Regular Expressions

- Why use meaning function?
  - Make clear what is syntax and what is semantics
  - Allow us to consider notation as a separate issue
    - Identifier:
      - Sequence of letters  $r_1$
      - Sequence of letters and digits  $r_2$

Integer: a non-empty sequence of digits

Integer: a non-empty sequence of digits

(0|1|2|3|4|5|6|7|8|9)



Integer: a non-empty sequence of digits

digit = (0|1|2|3|4|5|6|7|8|9)

Integer: a non-empty sequence of digits

digit = (0|1|2|3|4|5|6|7|8|9)

{digit}\*

Integer: a non-empty sequence of digits

digit = (0|1|2|3|4|5|6|7|8|9)

{digit}{digit}\*

Integer: a non-empty sequence of digits

digit = (0|1|2|3|4|5|6|7|8|9)

{digit}{digit}\*  
{digit}+

Identifier: sequence of letters or digits,  
starting with a letter

Identifier: sequence of letters or digits,  
starting with a letter

digit = (0|1|2|3|4|5|6|7|8|9)

Identifier: sequence of letters or digits,  
starting with a letter

digit = (0|1|2|3|4|5|6|7|8|9)

letter = (a|b|c|...|z|A|B|...|Z)

Identifier: sequence of letters or digits,  
starting with a letter

digit = (0|1|2|3|4|5|6|7|8|9)

letter = (a|b|c|...|z|A|B|...|Z)  
[a-z]



Identifier: sequence of letters or digits,  
starting with a letter

digit = (0|1|2|3|4|5|6|7|8|9)

letter = (a|b|c|...|z|A|B|...|Z)  
[a-zA-Z]

Identifier: sequence of letters or digits,  
starting with a letter

digit = [0-9]

letter = [a-zA-Z]

Identifier: sequence of letters or digits,  
starting with a letter

digit = [0-9]

letter = [a-zA-Z]

{letter}

Identifier: sequence of letters or digits,  
starting with a letter

digit = [0-9]

letter = [a-zA-Z]

$\{\text{letter}\}(\{\text{letter}\}|\{\text{digit}\})^*$

Whitespace: a non-empty sequence of  
blanks, newlines and tabs

Whitespace: a non-empty sequence of  
blanks, newlines and tabs

" "

Whitespace: a non-empty sequence of  
blanks, newlines and tabs

" " | "\t"

Whitespace: a non-empty sequence of  
blanks, newlines and tabs

" " | "\t" | "\n"



Whitespace: a non-empty sequence of  
blanks, newlines and tabs

`(" " | "\t" | "\n")+`

## Definition of Numbers

digit = [0-9]

digits = [0-9]+

opt\_frac = ("."{digits})?

opt\_exp = (E(\+|\-)?{digits})?

num = {digits}{opt\_frac}{opt\_exp}

## Definition of Numbers

digit = [0-9]

digits = [0-9]+

opt\_frac = ("."{digits})?

opt\_exp = (E(\+|\-){digits})?

num = {digits}{opt\_frac}{opt\_exp}

345 , 345.04 , 2.14+e7

# Lex regular expressions

Expression	Matches	Example	Using core operators
$c$	non-operator character $c$	$a$	
$\backslash c$	character $c$ literally	$\backslash *$	
$"s"$	string $s$ literally	$" * "$	
$.$	any character but newline	$a . * b$	
$^$	beginning of line	$^ abc$	used for matching
$\$$	end of line	$abc \$$	used for matching
$[s]$	any one of characters in string $s$	$[ abc ]$	$(a b c)$
$[^s]$	any one character not in string $s$	$[ ^ a ]$	$(b c)$ where $\square = \{a,b,c\}$
$r^*$	zero or more strings matching $r$	$a^*$	
$r^+$	one or more strings matching $r$	$a^+$	$aa^*$
$r^?$	zero or one $r$	$a^?$	$(a \epsilon)$
$r\{m,n\}$	between $m$ and $n$ occurrences of $r$	$a\{ 2 , 3 \}$	$(aa aaa)$
$r_1 r_2$	an $r_1$ followed by an $r_2$	$ab$	
$r_1   r_2$	an $r_1$ or an $r_2$	$a   b$	
$(r)$	same as $r$	$( a   b )$	
$r_1 / r_2$	$r_1$ when followed by an $r_2$	$abc / 123$	used for matching

# Regular Expressions for Lexical Analysis

- Write a regexp for the lexemes of each token class
  - Integer = `digit+`
  - Identifier = `letter(letter|digit)+`
  - Lparen = `'('`
  - ...
- Construct  $R$ , matching all lexemes for all tokens.
  - $R = \text{Integer} \mid \text{identifier} \mid \dots$   
 $= R_1 \mid R_2 \mid R_3 \mid \dots$

# Lexical Analyzer

1. Let input be  $x_1 x_2 \dots x_n$

For  $1 \leq i \leq n$  check

$$x_1 \dots x_i \in L(R)$$

$$R = R_1 \mid R_2 \mid \dots \mid R_n$$

2. If success, then we know that

$$x_1 \dots x_i \in L(R_j) \text{ for some } j$$

3. Remove  $x_1 \dots x_i$  from input and go to (1)

# Lexical Analyzer

- How much input is used ?

$x_1 \dots x_i \in L(R)$

$x_1 \dots x_j \in L(R)$   
 $i \neq j$

$\boxed{x_1 x_2 x_3 x_4} x_5 x_6 x_7 \dots$

**==**

**Maximal Munch**

$\max(i, j)$

**T\_ASSIGN**

**T\_EQ**



# Lexical Analyzer

- Which token is used? **Choose the one listed first!**

$$x_1 \dots x_i \in L(R)$$

$$x_1 \dots x_i \in L(R_j)$$

$$x_1 \dots x_i \in L(R_k)$$

$$R = R_1 \mid R_2 \mid \dots \mid R_n$$

IF = if

identifier = letter(letter | digit)\*

$$\text{if} \in \begin{matrix} L(\text{IF}) \\ L(\text{Identifier}) \end{matrix}$$

This situation arises when `if` is followed by `Lparen` e.g. `if (`

# Lexical Analyzer

- What if no rule matches?

$$x_1 \dots x_i \notin L(R)$$

Error = all strings not in specification

Put it last in priority

# Regexps in Lexical Analysis

- Regular expressions are a concise notation for string patterns
- Use in lexical analysis requires small extensions
  - Resolve ambiguities
  - Handle errors
- A good algorithm for lexical analysis will:
  - Require only single pass over the input
  - Few operations per character (lookup table)