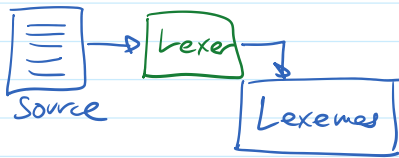


## 4 Syntax and Regular Expressions

Wednesday, January 31, 2024 11:57 AM



- From Linguistics:

- Orthography :- How to write words?

apple manzana pomme УЯЛЮ

given by a Dictionary:

In a programming language: keywords:

FOR WHILE FUNCTION.

- Syntax :- How to write sentences?

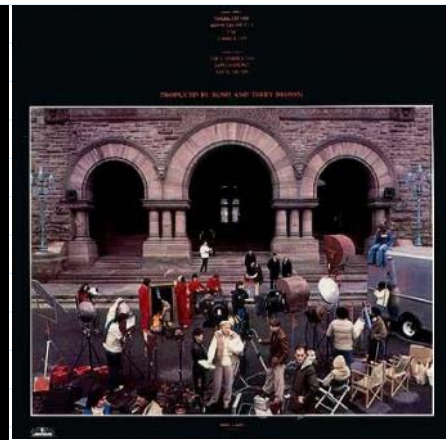
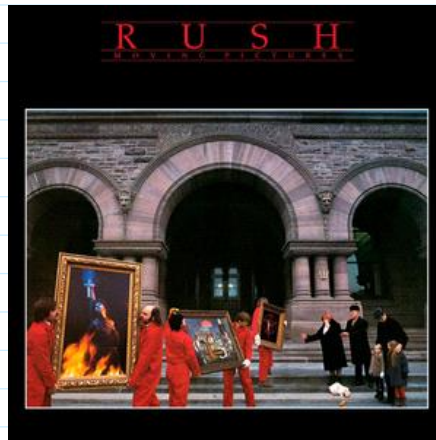
"has jump apple dog green"

- Semantics :- What do sentences mean?

"colorless green ideas sleep furiously" N. Chomsky  
Adj Adj noun verb adverb.

"the chickens are ready to eat"

"Moving Pictures"



• For a programming language.

- How to write words?

- Keywords :- for while if end ...

- How to write words?

- Keywords :- for while if end ...

- "identifiers" :- named entity.

- "Literals" :- numerical values  
strings  
{arrays}  
(dictionaries)  
(tuples)

integer  
reals  
hexes  
Scientific

To Answer this Question:

"Mathematical Languages"

S. Kleene  
E. Post.

Def. = Alphabet :- a finite set of symbols

e.g.  $\Sigma = \{a, b\}$

Language :- a set of sequences (words) made of symbols from the alphabet.

e.g.  $L = \{a, aa, aba, baa, bbb\}$

1. Specification Problem:

How to precisely describe a language?

$L = \{a, ab, aba, abb, \dots\}$   
???

2. Recognition Problem

given a word  $w$  and a language  $L$

answer  $w \in L$ ?

Answer 1 :- Specification :- Regular Expressions

built from:  $\Sigma$  | \* ( )

Regular Expressions:

- a word  $w$  is a regular expression
- if  $r_1$  and  $r_2$  are regular expressions
  - $r_1 r_2$
  - $r_1 | r_2$
  - $r_1^*$are also regular expressions

$- r_1 r_2$   
 $- r_1 | r_2$   
 $- r_1^*$

} are also regular expressions

A regular expression specifies a language

- $w$   $\{w\}$
  - $r_1 r_2$   $\{w_1 w_2 \mid w_1 \in r_1 \wedge w_2 \in r_2\}$
  - $r_1 | r_2$   $\{w \mid w \in r_1 \vee w \in r_2\}$
  - $r^*$   $\{w \mid w \text{ is zero or more repetitions of a word in } r\}$
- ↑ Kleene star.

Eg:  $(a | b) b = \{ab, bb\}$

$a : \{a\}$   
 $b : \{b\}$   
 $a|b : \{a, b\}$

Eg:  $a(b^*) = \{a, ab, abb, abbb, abbbb, \dots\}$

↑ specifying an infinite set through a finite expression.

Eg:  $a^* = \{\epsilon, a, aa, aaa, aaaa, aaaaa, \dots\}$

↑ empty string.

Eg:  $(aa | b^*) a = \{a, ba, aaa, bba, bbba, bbbba, \dots\}$

$\{aa\}$        $\{\epsilon, b, bb, \dots\}$

### • Shorthands

$r? \equiv (r | \epsilon)$

"to r or not to r"

$r+ \equiv r(r^*)$

"one or more repetitions of a word in r"

### • Regular Expressions AND Programming Languages

$\Sigma = \text{ASCII keyboard characters.}$

NOTE:- we cannot use regular expressions to describe a programming language.

$L_{C++} = \text{The set of all valid C++ programs.}$

Why?

consider  $\{a^n b^n \mid n \in \mathbb{Z}^+\} = \{ab, aabb, aaabbb, aaaa bbbb, \dots\}$

like matching brackets

$\{( ) [ ] \}$  note: counting is not enough.  
 $\{ ( [ ] ) \}$

Famous Stack Overflow Question:

"Regular Expression to validate HTML?"

**Closed**

$\langle a \rangle \dots \langle /a \rangle$   
 $\langle h3 \rangle \dots \langle /h3 \rangle$

• RegEx: are good enough to describe "tokens"

token:- a category of atomic strings.

lexeme:- an instance of a token.

- "identifiers" :- named entity.

- "Literals" :- numerical values  
strings.

integer  
reals  
hexes  
Scientific

while ( /flag/ == /62.3/ else if / ) /14/ 7 /! = x / + /var\_1/

Keyword    ident    op    real    literal    keyword    int    Lit    op    ident.

- Describing tokens using RegEx

$\Sigma = \text{keyboard characters}$

• Shorthands.

$[0-9] \equiv (0 | 1 | 2 | 3 | \dots | 8 | 9)$   
 $[A-Z] \equiv (A | B | C | \dots | Y | Z)$

$$\begin{aligned}
 [0-9] &\equiv (0|1|2|3|\dots|8|9) \\
 [A-Z] &\equiv (A|B|C|\dots|Y|Z) \\
 [a-z] &\equiv (a|b|c|\dots|y|z) \\
 [A-Za-z] &\equiv ([A-Z]|[a-z])
 \end{aligned}$$

- Regex for keywords.

e.g C++ (if | while | for | do | else | .....)

- Regex for integer literals

$$(+|-)?([0-9])^+ \left\{ \begin{array}{l} 0 \quad 1 \quad 123 \quad 001 \quad 000 \\ -5 \quad +27 \quad \quad -001 \quad +001 \end{array} \right\}$$

- Regex for identifiers.

Fortran  $([A-Z])([A-Z]|[0-9])^*$

A A2 M7 APPLE3

Pascal  $([A-Za-z])([A-Za-z]|[0-9]|' ')^*$

var\_1 point\_3D intX

Python  $([A-Za-z]|' ')([A-Za-z]|[0-9]|' ')^*$

var\_1 this\_is\_snake  
--init-- --main--

-----  
-----

- Regex for scientific notation.

e.g.  $3.7e^{-5}$   $-8.9e^7$   $123e^3$

$$(+|-)? [0-9]^+ (\cdot [0-9]^+)? 'e' (+|-)? [0-9]^+$$

- Regex can also be used to describe other standard formats.

- phone numbers
- SSN
- URLs
- Serial numbers
- licence plates
- etc.....

—●— EOF